# Wildfire Prediction

**Use Machine Learning to Predict Wildfire in Particular Area**



## Taniya Adhikari
**March 1, 2022**

# 1 INTRODUCTION

## 1.1 Objective

Identify the most important factors that cause wildfires and build a model to predict wildfire events accurately using MODIS fire products data, historical weather, and soil quality data.

## 1.2 Background

Every year, wildfires destroy massive number of homes and acres of lands and properties around the world. Wildfires are getting dangerously longer and severer each year. Wildfires have burned an average of 7 million acres yearly since 2000 in the U.S. Climate change has been a main factor in increasing wildfires in U.S. There are various reports from United Nations on rising temperatures across the globe are contributing to wildfires by increasing drought, decreasing precipitation, soil moisture, and the presence of trees and other potential fuel. By exploring these factors that cause wildfires, allows experts to take early action to minimize the risks and threats it poses. Though, climate change is the key factor, human related factors such as debris burning, campfire and arson is also a contributing to rising wildfires. It is estimated 80% of fires are caused by combination of these factors.

Wildfires are getting uncontrollable over the years, cause many fatalities and it is costly to contain it. In 2021, Dixie fire was one of the deadliest wildfires that happened in California area with 463,000 acres of land burned, which destroyed hundreds of buildings and threatened many lives. In 2020, Bay Area fire killed 35 people and burned almost 1 million acres of land. The history of wildfires goes all the way back to 1871 when Peshtigo Fire in Wisconsin burned 1.2 million acres and killed 1152 people.

In this project, I will be exploring MODIS Collection 6 Active Fire Data provided by NASA. MODIS fire products identify thermal anomalies, such as volcanoes or vegetation fires. It can identify 50% of fires accurately by capturing relevant fire pixels and its' accuracy increases with the size of fire. Thus, fire pixels are not necessarily true fires. Along with this data I will also be using historical daily climate summaries, other meteorological indicators also

provided by NASA LaRC. Lastly, I will also use Soil quality data. My goal is to create a model that can be used to do early fire detection using combination of variables from these datasets and detect true fires. For this project, I will be narrowing the research to just California State fire data.

## 1.3 Business Problem

Firefighting resource is crucial in containing wildfires and should be allocated correctly when needed. MODIS historical fire pixel data contains useful spatial information that can be used to identify patterns over time. Meteorological indicators can also provide information regarding weather conditions that causes wildfire. Soil quality indicators give information regarding drought possibilities. Though, MODIS can identify large fires with higher accuracy, it does have too many false alarms. It is expensive to use resources for false alarms. Thus, accurate prediction of fire is crucial to minimize the cost and for early preparation. For this project, I would like to answer following questions: **Do weather and soil conditions has any significance in wildfires? Can I identify true fires more accurately by combining weather and soil data to fire pixels data? Can I identify the fire prone area before fire happens, using historical fire pixels data combined with weather indicators and soil quality indicator?**

## 2   METHODS

### 2.1   Data Collection and Data Understanding

For this project, there was not any single existing dataset that had all the factors combined. I will be explaining the process of collecting different input variables for the model. I used multiple datasets to create a single dataset with as many factors as needed for this project. Following are the list of data:

**California Fire Geo Data**

To predict fire events, I had to map California fires correctly with weather, fire pixels and soil data for the given date. Furthermore, I must make sure that start dates of all fires are as

correct as possible. I used three different resources to cross-check start dates of the known fire events in California. Below is the list of datasets:

- California Fire Perimeters, shape files: https://gis.data.ca.gov/

- Wildland Fire Open data: https://data-nifc.opendata.arcgis.com/

- Cal Fire open data: https://www.fire.ca.gov/

Furthermore, I used geopandas library to merge Wildland Fire data (Point geometry) and Cal Fire data (Point geometry) to California Fire Perimeters (Polygon geometry) by finding nearest point to polygon. I used burned areas (in acres) as a threshold to make sure that fire points are not mapped too far. *(See Appendix A.1 for Fire events Data Preparation Code).* I identified 3803 fires that happened in last 10 years (2011-2020). Key variables in this dataset are **Fire Date and Fire location.**

Note: This is the most important step. The most challenging task was to get the fire dates corrected. None of the datasets had the accurate start date information, some fire events had dates that were way off than when it happened, and the start dates among all datasets were not the same for some fire events. To get this part corrected, I had to hard code some dates manually by doing lot of internet research such as finding news articles for that year and month and look for fire events using the "Name" variable. *(See Appendix A.1 for Fire Events Preparation Code).*

**MODIS Collection 6 Data**

MODIS Collection 6 Data comes as granules in hdf5 files. Though, I initially attempted to get the data by using NASA earth data. However, those files were large, and my device could not process it. Fortunately, Kaggle.com had complete MODIS Data from 2010-2020. Below is the list of MODIS data: *(Appendix A.2 MODIS Data Preparation).*

- Fire Pixel data: https://earthdata.nasa.gov/

- MODIS CSV file: https://www.kaggle.com/

After cleaning the date set, I identified 114,599 fire pixels for California Area between 2011-2020. Key variables for this dataset are:

| | |
|---|---|
| **Brightness T21 and T31** | Measure of photons at a wavelength received. measured in Kelvin |
| **FRP** | Measures the rate of radiant heat output from a fire. |
| **Scan and Track** | 1 km fire pixels but MODIS pixels get bigger toward the edge of scan. Scan and track are used to measure actual size of pixel |
| **Fire Pixel Latitude** | Centre of 1 km fire pixels but not the actual location. |
| **Date** | Date to accurately map it with fire events dataset. |
| **Confidence** | It ranges between 0 to 100%. It intends to help user to understand the quality of fire pixels. Whether they are high risk or low risk. They can be binned in to 3 classes: low, nominal or high |
| **Day or Night:** | Time of the day when the pixel was taken |
| **Hot Spot Type** | Area type: Vegetation, Volcanoes, Other Static Land Source and Offshore |
| **Month** | The fire event month |

**Soil and Meteorological Data**

Soil quality data and Meteorological data was collected from Kaggle.com as well. Soil and Meteorological datasets were big files, thus required a separate code. Altogether, there were 11,353,524 data points with 52 variables. After cleaning, combining and filtering for 10 year data I ended up with 277,628 data points with 25 variables. *(Appendix A.3: Soil and Meteorological Data Preparation).*

- Meteorological data: https://power.larc.nasa.gov/ , https://www.kaggle.com/

- Soil Dataset: https://www.fao.org/soils-portal/data-hub/soil-maps-and-databases/harmonized-world-soil-database-v12/en/

- Daily weather summary for California State Summary: https://www.noaa.gov/

This data was cleaned so there was not much to do except filtering the years and remove redundant variables, that had .80 or more correlation. Key variables are **Latitude, Longitude, elevation, slope, soil quality, oxygen availability, rain-fed cultivated land, precipitation, windspeed Humidity and temperature.**

## 2.2   Exploratory Data Analysis
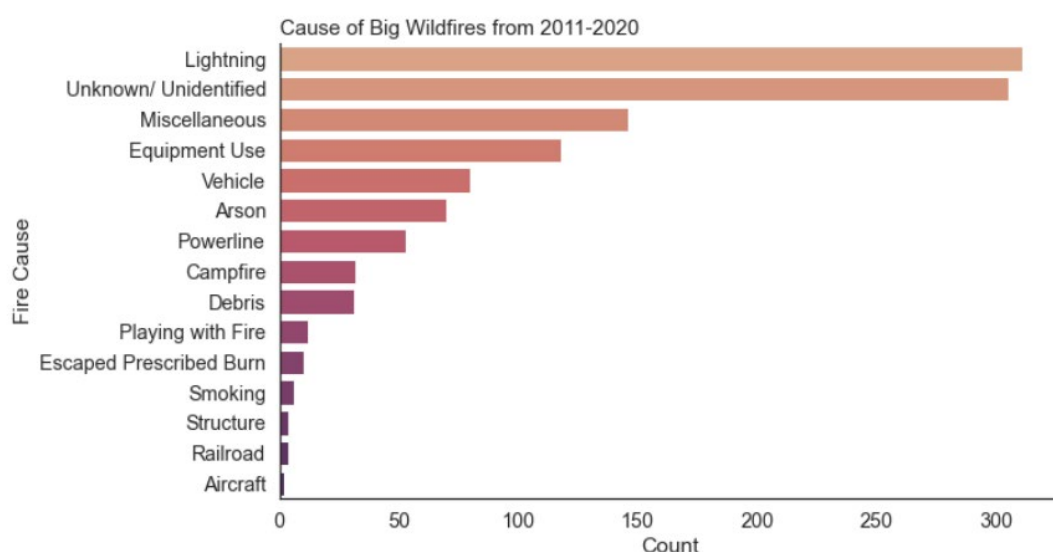
**Fire Events Data Analysis**

These are the key results I found:

There were **3803 fire events within 10 years of dataset for California** location. These are some of the important things I found when did a preliminary analysis:

**Figure 2.2.1 Causes of Big Wildfires Between 2011-2020**

**Lightning** was the number one reason to cause fire events that were bigger than 100 acres. *Note: 100 acres are considered small events and usually do not cause minimal harm.*
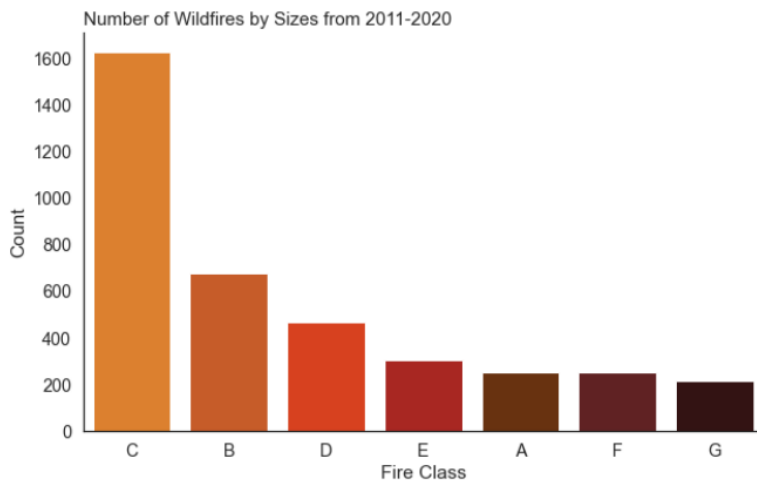
| Class A | 0 < area ≤ 0.25 Acres |
|---------|------------------------|
| Class B | 0.26 ≤ area ≤ 9.99 |
| Class C | 10 ≤ area ≤ 99.9 |
| Class D | 100 ≤ area ≤ 299 |
| Class E | 300 ≤ area ≤ 999 |
| Class F | 1000 ≤ area ≤ 4999 |
| Class G | 5000 ≤ |

Other reasons were Miscellaneous, equipment use, vehicle arson etc., indicating that human related behavior also contributes to big fires.

**Table 2.2.1** There are seven classes of fires based on area burned:

**Figure 2.2.2 Wildfire Occurrence by Class**

Class C wildfire had the highest occurrence during 2011-2020



**Figure 2.2.3 Top Reasons for Big Wildfires**

## Figure 2.2.4 Total Fires Between 2011-2020

2017 had the highest number of wildfires in past 10 years.



Total Fires in California (2011-2020)

## Figure 2.2.5 Wildfires Occurrence by Months

July month had the highest number of wildfires



Total Fires for Each Month (2011-2020)

**Figure 2.2.6: California Fires by Region.**



## Modis Collection 6 Data Analysis

These were the following analysis results:

For numerical data, Scan and track had the highest correlation 0f .99, so I took that variable out. Although Brightness T21 and T31 seems to be the same thing, but they have weak relationship, so I kept both variables.

**Figure 2.2.7 Distribution of Fire Pixels Attributes**

Distribution of Fire Pixel Attributes in West Coast Region (2011-2020)

## Figure 2.2.8 Fir Pixels Mapped



Geospatial Plot Calfornia Fire Pixels (2011-2020)

## Figure 2.2.9 Fire Pixels by Type of Hot Spot Area

Fire Pixels data indicates Vegetation area is more prone to high-risk wildfires



Fire Pixels data shows most fire pixels were taken in August, as opposed to **July fire months in Figure 2.2.5**

**Soil and Meteorological Data Analysis**

Here the few things I found. For Soil, data it was around the same region where most fire happened.

**Figure 2.2.11 Geospatial Plot of Temporal Soil Data**



For Meteorological Data. Here are the key findings:

**Figure 2.2.12 Comparison of Precipitation, Temperature and Windspeed.**

**June, July and August is when I have highest temperature and lowest precipitation,** in Figure 2.2.6, I found that Most fires occur in June, July and August, with July being the month with highest fire occurance.

# Figure 2.2.13 Comparison of Surface Pressure, Earth Skin Temperature, and

These also shows similar indication as the Figure 2.2.12



Metereological Indicators: SP vs. TS vs. WS Range



Average Monthly: Precipitation vs. Fires



Average Monthly: Temperatures vs. Fires

Average Monthly: WindSpeed vs. Fires

**Figure 2.2.15: Heat Map to show correlation**



## 2.3   Data Merging

To prepare data for analysis and predict fire events correctly, I must first make sure that all fire dates are as correct as possible. California Fire Perimeter dataset is a shapefile, which had fire area as polygon with latitude and longitude coordinates of each point. All the above datasets came with coordinates information (latitude and longitude). After cleaning the fires data (handing missing values, fixing dates, identifying unique fires), I narrowed all the data for

only California coordinates and filtered datasets for fires only between year 2011-2020. Furthermore, I converted fire pixels, weather, and soil data into geospatial data frame using `Geographic 2D CRS: EPSG:4326`. Then, I converted all of them including fire data into `Projected CRS: EPSG:3310,` so I can merge them by nearest point location (Note: This is crucial to get the minimum distance correctly in m or km, otherwise distance in degree unit holds no meaning. EPSG:3310 converts geographic coordinates into meters).

Steps into Data Preparation

1.  I combined daily temperatures data to MODIS pixels data by finding the nearest geo point on the given date of the pixel.

2. I combined soil and meteorological data the same way (all these data were daily data, thus merging them was not as complex).

3. After merging the above datasets, I merged fire data with combined dataset based on fire pixels geo point to nearest fire polygon for the same year and month (Note: I did not merge by exact date, because fire pixels can also be true fire event after it already happened. Fires can go on for weeks and sometimes months depending on the fire size). After merging all the data, I performed few other analyses and made few assumptions to make sure fires are mapped correctly *(See Section 2.7 Assumptions)*.

After merging the data, I found 113912 fire pixels were mapped and 687 was un-mapped. I ended up with total 63 input variables. To find out which fires were accurately mapped, I did separate analysis using assumptions from Section 10.1 Fires that were mapped outside of its burned area distance, were marked as false alarms. All the fires that were mapped correctly was considered labelled data and all the ones that weren't marked correctly were considered as un-labelled data. For Data Merging *See Appendix A.4 Data Merging Code.*

## 2.4  Data Modeling

Total data points were: 45186 labelled and 57935 un-labelled data and for California. Below is the plot of True and False Fire Pixels from MODIS Data.

**Figure 1.7.1 Geospatial analysis of labelled and un-labelled MODIS fire pixels data.**



**Figure 1.7.2 Geospatial analysis of Actual Fire Area and True Mapped Fire pixels.**

Using above visualization, it is safe to assume that that fire pixels are mapped correctly to fire events. Before data modeling step, one last thing I did is check for multi-collinearity since, some variables were duplicated due to data collection from multiple sources. There was total **25 variables.**

**Preliminary Data Modeling**

For Data modeling, I needed two-part analysis, because I had two things I wanted to answer, can I identify true fire pixels using current data? and can I do early prediction of wildfire using previous days data?

For the first part of the problem statement requires simple classification, I split the data into train and test set using sklearn package and stratified the data, so I don't have unbalanced classes. 70% of data were in train set and 30 % in test. For the binary classification, I chose two algorithms: **Random Forest Classifier and Support Vector Machines with RBF kernel (Non-Linear data model).** For SVM I also did one-hot-coding for categorical data and scaled input variables. *(See Appendix A.5 Data Modeling Part 1).* I also did cross validation to check for model accuracy.

*Note: In final analysis I will be doing model selection between these two models*

For the second part, I needed to do **Time Series Classification**. For time series classification, idea was to use last 3-7 days data to do early fire detection. Whether the fire will happen or not? I chose **LSTM RNN Model** for this approach. I did some feature engineering to transform the data, by shifting and prepending data from previous days, turned data into 3 dimensions. I split the data into train, test, and validation set. and build the model with 4 layers. *(See Appendix A.6 Data Modeling Part 2)*

*Note: In final Analysis I will be doing feature selection to see if the accuracy score increase.*

## 2.5 Preliminary Analysis

```
Accuracy score: 99.23%
Cross validation Accuracy score: 98.81%
Cross validation Precision score: 92.69%
Cross validation Recall score: 37.87%
Cross validation F1 score: 53.74%
```

Confusion Matrix with labels



These are the results from Random Forest Classification. Random Forest shows 99.23% accuracy with high precision score of 92.86%, recall score is 37.87% which is low and F1 score is 53.74%

```
Accuracy score: 98.3%
Cross validation Accuracy score: 98.26%
Cross validation Precision score: 86.29%
Cross validation Recall score: 5.43%
Cross validation F1 score: 10.18%
```

Confusion Matrix with labels



For SVM model Accuracy score was slight lower than Random Forest 98.3%, precision is 86.29%, recall score is 5.43% and F1 score 10.18%

**For Time Series here are the results**

I did analysis with 3 days, 5 days, and 7 days data.

### 3 days LSTM model

**Before Feature Selection**

```
Accuracy score: 64.36%
Precision score: 22.19%
Recall score: 16.74%
```

### 5 days LSTM Model

**Before Feature Selection**

```
Accuracy score: 67.08%
Precision score: 28.86%
Recall score: 21.23%
```

### 7 days LSTM Model

**Before Feature Selection**

```
Accuracy score: 64.6%
Precision score: 21.7%
Recall score: 15.71%
```

Testing score **using 5 days data shows a slightly higher accuracy with 67.08% and precision 28.86%.** I only attached results form 5 days because that showed highest score. I also run the model for 3 days and 7 days data, but the accuracy score was too low. *See Appendix A.6 Data Modelling Part 2 for full code and results*

## 2.6   Conclusion

**Do weather and soil conditions have any significance in wildfires?**

Yes, in Figure **2.2.14** Ihighest temperature months were between Jun, July, and August when the highest number of fires happened. Not only that precipitation also had negative relationship to number of wildfires. Not having enough precipitation can contribute to drought which can result into wildfires.

**Can I identify true fires more accurately by combining weather and soil data to fire pixels data?**

Yes, as of now MODIS accuracy is about 50% and its' prediction gets accurate as the fire becomes bigger. Whereas the random forest model increases the accuracy when combined with temperature and soil data. However, I would like to mention recall score was below 50% meaning there are high number of false negative.

**Can I identify the fire prone area before fire happens, using historical fire pixels data combined with weather indicators and soil quality indicator?**

Maybe, this is the problem of time series classification I use previous 5 days to see if I can predict whether it is a true fire or not. Model shows 67.08% accuracy. It is not as high but better than 50 %. However, I have not done feature selection yet, for the final paper, I will be including those results. *(Note: Please see the final paper for results)*

## 2.7   Assumptions

MODIS data collects pixels and its' location. Each location is 1 km away from the centroid of the pixel. It is never exact location of the fire, from what I know there could be two fires and it could belong one or the other. In addition, MODIS also detects smoke, if there are big wildfires, MODIS can detect a pixel as far as the smoke is going. For simplicity when I mapped wildfires with fire pixels, I mapped them by nearest point to polygon merge. I assumed that any fire that is mapped and is within their burned area, I marked it as labelled and true fire data. For example, if the fire is 10 km away and total burned area is 10 km sq (2417 acres) then I marked it as true fire pixel.

## 2.8   Limitations and Challenges

There were many Challenges and Limitations:

1.  Data Collection took significant amount of time. Downloading data from NASA database and interpreting data was difficult and time consuming, ultimately, I ended up choosing easier route and found a database that was already clean.

2.  Not having the exact date of fire event or having dates that are not correctly or missing dates that had end of the year date as its start date. Finding and cleaning

that data took a huge amount of time and work, even though I got the data from official state websites. Data was not managed and had too many missing values.

3. Fire size was measured using acres or km square, however fire distance was 1 dimensional, so mapping it without assumption was difficult.

4. Some fires were mapped 3000-4000 km when the fire size was only 3-4 km. Identifying them and marking them as un-labelled data took significant amount of time.

## 2.9   Future Applications

For future applications, I do believe the model can be better with the right data and clean data. The most important thing is to have correct date fore fire. I also think including data from other states can make the performance better.

## 2.10 Recommendations

I mapped fires by year and month, but to further build the model I believe fire should be mapped exactly on the date to do early prediction. Perhaps, LSTM was not a right approach, there is a random Forest time series classification model, that may perform better.

## 2.11 Implementation

This can be implemented along with MODIS to accurately predict true fire pixels, so fire fighting resources can be allocated correctly. As mentioned earlier MODIS detects 50% of the fires accurately, and only large fires are accurately detected, but adding external factors such as weather and soil quality can increase the accuracy,

## 2.12 Ethical Assessment

Because wildfires are dangerous and can be a threat to human life. It is crucial to choose right data and correct modelling approach. For this project, I collected data from credible resources such as NASA Earth Data. For modelling.

## 3   REFERENCES

1.    *15 worst wildfires in US history: Earth.org - past: Present: Future*. Earth.Org. (2021,
      August 30). Retrieved February 16, 2022, from https://earth.org/worst-wildfires-in-us-
      history/

2.    Ares. (2020, February 9). *California wildfires (2013-2020)*. Kaggle. Retrieved
      February 17,  2022, from https://www.kaggle.com/ananthu017/california-wildfire-
      incidents-20132020

3.    California Department of Forestry and Fire Protection (CAL FIRE). (n.d.). *Incidents
      overview*. Cal Fire Department of Forestry and Fire Protection. Retrieved February
      17, 2022, from https://www.fire.ca.gov/incidents/

4.    *California fire perimeters (all)*. California State Geoportal. (n.d.). Retrieved February
      17, 2022, from https://gis.data.ca.gov/maps/CALFIRE-Forestry::california-fire-
      perimeters-all-1/about

5.    Chen, F., Niu, S., Tong, X., Zhao, J., Sun, Y., & He, T. (2014, August 27). *The
      impact of precipitation regimes on forest fires in Yunnan Province, Southwest China*.
      The Scientific World Journal. Retrieved February 17, 2022, from
      https://www.hindawi.com/journals/tswj/2014/326782/

6.    Eswaramoorthy, V. D. (n.d.). *Increasing the stStatistical Significance For Modis
      Active Fire Hotspots in Portugal Using One-Class Support Vector Machines*.
      Retrieved February 17, 2022, from https://cartographymaster.eu/wp-
      content/theses/2017_Eswaramoorthy_Thesis.pdf

7.    *Harmonized World Soil Database v 1.2*. FAO SOILS PORTAL. (n.d.). Retrieved
      February 17, 2022, from https://www.fao.org/soils-portal/data-hub/soil-maps-and-
      databases/harmonized-world-soil-database-v12/en/

8.      NASA. (2022, February 17). *MCD14DL*. NASA. Retrieved February 17, 2022, from

https://earthdata.nasa.gov/earth-observation-data/near-real-time/firms/c6-

mcd14dl#ed-firms-attributes

9.      NASA. (2022, February 17). *Wildfires*. NASA. Retrieved February 17, 2022, from

https://earthdata.nasa.gov/learn/toolkits/wildfires

10.     NASA. (n.d.). *NASA Power*. NASA. Retrieved February 17, 2022, from

https://power.larc.nasa.gov/

11.     *The national burn rate is going up*. Face the Facts USA. (2013, August 26). Retrieved

February 17, 2022, from https://facethefactsusa.org/facts/the-national-burn-rate-is-

going-up--literally/

12.     National Centers for Environmental Information (NCEI). (n.d.). *Climate Data Online

Search*. Search | Climate Data Online (CDO) | National Climatic Data Center

(NCDC). Retrieved February 17, 2022, from https://www.ncdc.noaa.gov/cdo-

web/search

13.     *WFIGS - Wildland Fire Locations Full History*. National Interagency Fire Center.

(n.d.). Retrieved February 17, 2022, from https://data-

nifc.opendata.arcgis.com/datasets/wfigs-wildland-fire-locations-full-history/explore

14.     *Wildfires and climate change*. Center for Climate and Energy Solutions. (2021, July

22). Retrieved February 17, 2022, from https://www.c2es.org/content/wildfires-and-

climate-change/

15.     Ye, A. (2020, October 3). *NASA Wildfire Satellite Data*. Kaggle. Retrieved February

17, 2022, from https://www.kaggle.com/washingtongold/wildfire-satellite-data

16.     Zaitseff, S. (2019, January 18). *Classification of time series with LSTM RNN*. Kaggle.

        Retrieved February 17, 2022, from https://www.kaggle.com/szaitseff/classification-

        of-time-series-with-lstm-rnn

## 4   APPENDICES

Please see the attached documents

A.1 Fire events Data Preparation

A.2 MODIS Collection 6 Data Preparation

A.3 Soil and Meteorological Data Preparation

A.4 Data Merging Code

A.5 Data Modelling Part 1 – Simple Classification

A.6 Data Modelling Part 2 – Time Series Classification

# Appendix A.1 Fire Events Data Preparation code

Following code provides the steps that were taken to prepare the fire events data with correct dates and location. We identified 3803 wildfire events happened between 2011-2020. Out of which 779 fires were bigger than 300 acres.

```python
In [1]: import datetime as dt
        from pathlib import Path
        import math
        import os
        import sqlite3
        import json
        import geopandas as gpd
        import pygeos
        import pyproj
        import shapely
        import shapely.ops as ops
        from shapely.geometry import Point, Polygon
        from shapely.geometry.polygon import Polygon
        from functools import partial

        import pandas as pd
        import numpy as np
        import seaborn as sns
        import matplotlib.pyplot as plt
        %matplotlib inline

        from sklearn.model_selection import train_test_split


        from sklearn import svm
        from sklearn.svm import SVC
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.naive_bayes import GaussianNB
        from sklearn.metrics import accuracy_score, classification_report, confusion_m
        atrix

        from sklearn.feature_selection import SelectKBest
        from sklearn.feature_selection import chi2, f_classif, mutual_info_classif
        from functools import partial


        from sklearn.preprocessing import StandardScaler

        import warnings
        warnings.filterwarnings('ignore')
```

# Data Collection of all 3 Fire events Dataset

**1a. Load the dataset 1: All California Fire Dataset - https://gis.data.ca.gov/ (https://gis.data.ca.gov/)**

```
In [2]: ca_fires_df = gpd.read_file("Data/California_Fire_Perimeters_(all).geojson")
        ca_fires_df.head(2)
```

Out[2]:

| | OBJECTID | YEAR_ | STATE | AGENCY | UNIT_ID | FIRE_NAME | INC_NUM | ALARM_DATE | |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 21440 | 2020 | CA | CDF | NEU | NELSON | 00013212 | 2020-06-18T00:00:00+00:00 | 23 |
| **1** | 21441 | 2020 | CA | CDF | NEU | AMORUSO | 00011799 | 2020-06-01T00:00:00+00:00 | 04 |

```
In [3]: ca_fires_df.shape
```

Out[3]: (21318, 19)

```
In [4]: ca_fires_df.crs
```

```
Out[4]: <Geographic 2D CRS: EPSG:4326>
        Name: WGS 84
        Axis Info [ellipsoidal]:
        - Lat[north]: Geodetic latitude (degree)
        - Lon[east]: Geodetic longitude (degree)
        Area of Use:
        - name: World.
        - bounds: (-180.0, -90.0, 180.0, 90.0)
        Datum: World Geodetic System 1984 ensemble
        - Ellipsoid: WGS 84
        - Prime Meridian: Greenwich
```

**1b. Load the dataset 1: Supplementary Data for Fire Incident Location - https://data-nifc.opendata.arcgis.com/ (https://data-nifc.opendata.arcgis.com/)**

```
In [5]: fire_location = gpd.read_file("Data/WFIGS_-_Wildland_Fire_Locations_Full_Histo
        ry.geojson")
        print(fire_location.shape)
        fire_location.head(2)
```

(208013, 95)

Out[5]:

| | OBJECTID | ABCDMisc | ADSPermissionState | CalculatedAcres | ContainmentDateTime | Controll |
|---|---|---|---|---|---|---|
| **0** | 1 | None | CERTIFIED | 50.64 | 2020-08-06T23:13:07+00:00 | 06T23:13: |
| **1** | 2 | None | DEFAULT | NaN | None | |

2 rows × 95 columns

**1c. Load the dataset 1: Supplementary California Fire Dataset - https://www.fire.ca.gov/ (https://www.fire.ca.gov/)**

```
In [6]: fires_df = pd.read_csv('Data/California_Fire_Incidents.csv')
        print(fires_df.shape)
```

(1636, 40)

**USA Shape File https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html (https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html)**

```
In [7]:  USA = gpd.read_file("Data/County/cb_2018_us_county_500k.shp")
         USA.head()
```

Out[7]:

| | STATEFP | COUNTYFP | COUNTYNS | AFFGEOID | GEOID | NAME | LSAD | ALAND | AV |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 21 | 007 | 00516850 | 0500000US21007 | 21007 | Ballard | 06 | 639387454 | 694 |
| 1 | 21 | 017 | 00516855 | 0500000US21017 | 21017 | Bourbon | 06 | 750439351 | 48 |
| 2 | 21 | 031 | 00516862 | 0500000US21031 | 21031 | Butler | 06 | 1103571974 | 139 |
| 3 | 21 | 065 | 00516879 | 0500000US21065 | 21065 | Estill | 06 | 655509930 | 65 |
| 4 | 21 | 069 | 00516881 | 0500000US21069 | 21069 | Fleming | 06 | 902727151 | 7 |

## Data Preliminary Analysis

```
In [8]:  # check for missing value
         def percentMissing(df):

             df_numeric = df.select_dtypes(include=[np.number])
             numeric_cols = df_numeric.columns.values

             # % of missing data
             for col in df.columns:
                 # create missing indicator for features with missing data
                 missing = df[col].isnull()
                 pct_missing = np.mean(missing)*100
                 #if pct_missing >60:
                 print('{} - {}%'.format(col, round(pct_missing)))
                 num_missing = np.sum(missing)
```

```
In [9]:   # Checking data type
          def Datatype(df):
              # shape and data types of the data
              print("There are {} rows and {} columns".format(df.shape[0], df.shape[1]))
              print(df.dtypes)

              # select numeric columns
              df_numeric = df.select_dtypes(include=[np.number])
              numeric_cols = df_numeric.columns.values
              print(numeric_cols)

              # select non numeric columns
              df_non_numeric = df.select_dtypes(exclude=[np.number])
              non_numeric_cols = df_non_numeric.columns.values
              print(non_numeric_cols)
```

# Data Exploration: All California Fire Dataset

## All Fire Incident Geopoint data to map fire with right date and time for the first discovery or reported fire incident

```
In [10]:  fire_location = fire_location[(fire_location['InitialLatitude']<= 42) & (fire_
          location ['InitialLatitude'] >= 32)]
          fire_location = fire_location[(fire_location['InitialLongitude']<= -114) & (fi
          re_location['InitialLongitude'] >= -126)]
```

```
In [11]:  fire_location = fire_location[['InitialLatitude','InitialLongitude','FireDisco
          veryDateTime','ContainmentDateTime',
                                          'ControlDateTime','FireOutDateTime','POOState',
          'FireCause','GACC','IncidentName',
                                          'LocalIncidentIdentifier','UniqueFireIdentifie
          r','WFDSSDecisionStatus','geometry']]
```

```
In [12]:  import string
          import re
          def get_year(text):
              pattern = r'[^A-Za-z ]'
              if re.search("[^0-9]", text):
                  year = text[0:4]
              else:
                  None
              return year
```

```
In [13]:  fire_location['IncidentYear'] = fire_location['UniqueFireIdentifier'].apply(la
          mbda x: get_year(x))
```

```
In [14]:  fire_location['IncidentYear'] = fire_location['IncidentYear'].astype(int)
```

```
In [15]: fire_location = fire_location[(fire_location['IncidentYear'] >=2010) & (fire_l
         ocation['POOState']=='US-CA')]
         fire_location = fire_location[fire_location['IncidentYear'] <=2020]
```

```
In [16]: fire_location = fire_location.drop(['ContainmentDateTime', 'ControlDateTime',
                                             'FireOutDateTime', 'WFDSSDecisionStatus'],
         axis = 1)
```

```
In [17]: fire_location.crs
```

```
Out[17]: <Geographic 2D CRS: EPSG:4326>
         Name: WGS 84
         Axis Info [ellipsoidal]:
         - Lat[north]: Geodetic latitude (degree)
         - Lon[east]: Geodetic longitude (degree)
         Area of Use:
         - name: World.
         - bounds: (-180.0, -90.0, 180.0, 90.0)
         Datum: World Geodetic System 1984 ensemble
         - Ellipsoid: WGS 84
         - Prime Meridian: Greenwich
```

```
In [18]: fire_location['DiscoveryDate'] = fire_location['FireDiscoveryDateTime'].astype
         ('datetime64[ns]')
         fire_location['DiscoveryDate']= fire_location['DiscoveryDate'].dt.strftime('%Y
         -%m-%d')
         fire_location['DiscoveryDate'] = fire_location['DiscoveryDate'].astype('dateti
         me64[ns]')
         fire_location['DiscoveryYear'] = fire_location['DiscoveryDate'].dt.year
         fire_location['DiscoveryMonth'] = fire_location['DiscoveryDate'].dt.month
         fire_location['DiscoveryDay'] = fire_location['DiscoveryDate'].dt.day
```

```
In [19]: fire_location = fire_location[['UniqueFireIdentifier', 'IncidentYear', 'Incide
         ntName', 'DiscoveryDate',
                                         'DiscoveryYear','DiscoveryMonth', 'DiscoveryDa
         y', 'geometry']]
```

```
In [22]: fire_location.head(2)
```

Out[22]:

|  | UniqueFireIdentifier | IncidentYear | IncidentName | DiscoveryDate | DiscoveryYear | Discovery |
|---|---|---|---|---|---|---|
| **79385** | 2011-NVCNC-000020 | 2011 | Washoe | 2011-01-26 | 2011 | |
| **128205** | 2014-CALBOR-001660 | 2014 | Casitas | 2014-06-17 | 2014 | |

```
In [23]: fire_location = fire_location.sort_values(['DiscoveryDate'], ascending=True)
```

```
In [24]: percentMissing(fire_location)
```

```
UniqueFireIdentifier - 0%
IncidentYear - 0%
IncidentName - 0%
DiscoveryDate - 0%
DiscoveryYear - 0%
DiscoveryMonth - 0%
DiscoveryDay - 0%
geometry - 0%
```

Check for Correct year for fire incidents and delete duplicates based on coordinates and date

```
In [25]: fire_location["Disc_minus_ID"] = fire_location["IncidentYear"] - fire_location
         ["DiscoveryYear"]
         fire_location[fire_location["Disc_minus_ID"]!=0]
```

Out[25]:

| | UniqueFireIdentifier | IncidentYear | IncidentName | DiscoveryDate | DiscoveryYear | DiscoveryI |
|---|---|---|---|---|---|---|
| **53158** | 2016-CACNF-002667 | 2016 | MUTUAL AID DE LUZ | 2015-08-09 | 2015 | |

```
In [26]: fire_location = fire_location.drop(['Disc_minus_ID'], axis = 1)
```

```
In [27]: fire_location= fire_location[~fire_location.duplicated(['geometry', 'Discovery
         Date'], keep='first')]
```

```
In [28]: fire_location['DiscoveryYear'].describe()
```

```
Out[28]: count    23849.000000
         mean      2018.729129
         std          1.439810
         min       2011.000000
         25%       2018.000000
         50%       2019.000000
         75%       2020.000000
         max       2020.000000
         Name: DiscoveryYear, dtype: float64
```

# All Geospatial data for fire area perimeter

In [29]:
```
Datatype(ca_fires_df)
```

```
There are 21318 rows and 19 columns
OBJECTID              int64
YEAR_                object
STATE                object
AGENCY               object
UNIT_ID              object
FIRE_NAME            object
INC_NUM              object
ALARM_DATE           object
CONT_DATE            object
CAUSE               float64
COMMENTS             object
REPORT_AC           float64
GIS_ACRES           float64
C_METHOD            float64
OBJECTIVE           float64
FIRE_NUM             object
SHAPE_Length        float64
SHAPE_Area          float64
geometry           geometry
dtype: object
['OBJECTID' 'CAUSE' 'REPORT_AC' 'GIS_ACRES' 'C_METHOD' 'OBJECTIVE'
 'SHAPE_Length' 'SHAPE_Area']
['YEAR_' 'STATE' 'AGENCY' 'UNIT_ID' 'FIRE_NAME' 'INC_NUM' 'ALARM_DATE'
 'CONT_DATE' 'COMMENTS' 'FIRE_NUM' 'geometry']
```

In [30]:
```python
ca_fires_df['ALARM_DATE'] = pd.to_datetime(ca_fires_df['ALARM_DATE'], errors = 'coerce')
ca_fires_df['ALARM_DATE'] = ca_fires_df['ALARM_DATE'].astype('datetime64[ns]')
ca_fires_df['FireDate']= ca_fires_df['ALARM_DATE'].dt.strftime('%Y-%m-%d')
ca_fires_df['FireDate'] = ca_fires_df['FireDate'].astype('datetime64[ns]')
ca_fires_df['FireYear'] = ca_fires_df['FireDate'].dt.year
ca_fires_df['FireMonth'] = ca_fires_df['FireDate'].dt.month
ca_fires_df['FireDay'] = ca_fires_df['FireDate'].dt.day

ca_fires_df['CONT_DATE'] = pd.to_datetime(ca_fires_df['CONT_DATE'], errors = 'coerce')
ca_fires_df['ContDate']= ca_fires_df['CONT_DATE'].dt.strftime('%Y-%m-%d')
```

In [31]:
```python
ca_fires_df = ca_fires_df[(ca_fires_df['FireYear']>= 2011) & (ca_fires_df['STATE']== 'CA')]
ca_fires_df.shape
```

Out[31]:
```
(3677, 24)
```

In [32]:
```python
ca_fires_df = ca_fires_df.sort_values('FireDate')
```

```
In [33]:   # These columns are unnecessary
           # information, agency who was incharge is not needed,
           # INC NUM is unique, C_method is how it was tracked.

           ca_fires_df = ca_fires_df.drop(['YEAR_', 'C_METHOD', 'AGENCY', 'INC_NUM', 'ALA
           RM_DATE',
                                            'CONT_DATE','FIRE_NUM', 'COMMENTS', 'OBJECTIV
           E', 'SHAPE_Area',
                                            'SHAPE_Length', 'REPORT_AC'], axis = 1)
```

```
In [34]:   ca_fires_df = ca_fires_df.rename(columns={'OBJECTID': 'ObjectID', 'CAUSE': 'Fi
           reCause',
                                                       'GIS_ACRES': 'TotalAcres','STATE':'S
           tate', 'UNIT_ID':'UnitID',
                                                       'FIRE_NAME':'Name', 'ContDate': 'Con
           tainmentDate'})
```

```
In [35]:   percentMissing(ca_fires_df)

           ObjectID - 0%
           State - 0%
           UnitID - 1%
           Name - 0%
           FireCause - 1%
           TotalAcres - 0%
           geometry - 0%
           FireDate - 0%
           FireYear - 0%
           FireMonth - 0%
           FireDay - 0%
           ContainmentDate - 1%
```

Check for duplicate values based on geometry and date

```
In [36]:   ca_fires_df = ca_fires_df.sort_values(['FireDate'], ascending=True)
```

```
In [37]:   ca_fires_df = ca_fires_df[~ca_fires_df.duplicated(['geometry', 'FireDate'], ke
           ep='first')]
```

```
In [38]:   ca_fires_df = ca_fires_df[~ca_fires_df.duplicated(['geometry', 'TotalAcres'],
           keep='first')]
```

```
In [39]:   ca_fires_df.shape
```

```
Out[39]:   (3673, 12)
```

## All California Incidents with name and coordinates. This data is used as supplement data to check for fire date inaccuracies

In [40]:
```python
fires_df['Started'] = fires_df['Started'].astype('datetime64[ns]')
fires_df['CaDate']= fires_df['Started'].dt.strftime('%Y-%m-%d')
fires_df['Extinguished'] = fires_df['Extinguished'].astype('datetime64[ns]')
fires_df['ExitDate']= fires_df['Extinguished'].dt.strftime('%Y-%m-%d')
```

In [41]:
```python
fires_df = fires_df[(fires_df['Latitude']<= 42) & (fires_df['Latitude'] >= 32
)]
fires_df = fires_df[(fires_df['Longitude']<= -114) & (fires_df['Longitude'] >=
-126)]
```

In [42]:
```python
fires_df['CaDate'] = fires_df['CaDate'].astype('datetime64[ns]')
fires_df['CaYear'] = fires_df['CaDate'].dt.year
fires_df['CaMonth'] = fires_df['CaDate'].dt.month
fires_df['CaDay'] = fires_df['CaDate'].dt.day
```

In [43]:
```python
fires_df = fires_df[['Name', 'Latitude', 'Longitude', 'CaDate','CaYear', 'CaMo
nth', 'CaDay', 'ExitDate', 'AcresBurned',
                     'ArchiveYear', 'Counties', 'UniqueId']]
fires_df.shape
```

Out[43]: (1465, 12)

In [44]:
```python
percentMissing(fires_df)
```

```
Name - 0%
Latitude - 0%
Longitude - 0%
CaDate - 0%
CaYear - 0%
CaMonth - 0%
CaDay - 0%
ExitDate - 4%
AcresBurned - 0%
ArchiveYear - 0%
Counties - 0%
UniqueId - 0%
```

Check for Correct year for fire incidents and delete duplicates based on coordinates and date

In [45]:
```python
fires_df["Archive_minus_Year"] = fires_df["ArchiveYear"] - fires_df["CaYear"]
fires_df[fires_df["Archive_minus_Year"]!=0]
```

Out[45]:

| | Name | Latitude | Longitude | CaDate | CaYear | CaMonth | CaDay | ExitDate | AcresBurned | A |
|---|---|---|---|---|---|---|---|---|---|---|
| **1019** | Taglio Fire | 37.21812 | -121.07761 | 1969-12-31 | 1969 | 12 | 31 | 2018-01-09 | 12.0 | |
| **1261** | Bridge Fire | 38.07135 | -122.76751 | 1969-12-31 | 1969 | 12 | 31 | 2019-01-04 | 45.0 | |

In [46]:
```python
idx = fires_df[fires_df["Archive_minus_Year"]!=0].index.tolist()

# Taglio fire was in May 17, 2017
fires_df.at[idx[0],'CaDate']= '2017-05-17'
fires_df.at[idx[0],'ExitDate']= '2017-05-17'
fires_df.at[idx[0],'CaYear']= 2017
fires_df.at[idx[0],'CaMonth']= 5
fires_df.at[idx[0],'CaDay']= 17


# Taglio fire was in May 17, 2017
fires_df.at[idx[1],'CaDate']= '2018-08-08'
fires_df.at[idx[1],'ExitDate']= '2018-08-09'
fires_df.at[idx[1],'CaYear']= 2018
fires_df.at[idx[1],'CaMonth']= 8
fires_df.at[idx[1],'CaDay']= 8
```

In [47]:
```python
fires_df["Archive_minus_Year"] = fires_df["ArchiveYear"] - fires_df["CaYear"]
fires_df[fires_df["Archive_minus_Year"]!=0]
```

Out[47]:

| Name | Latitude | Longitude | CaDate | CaYear | CaMonth | CaDay | ExitDate | AcresBurned | Archive |
|---|---|---|---|---|---|---|---|---|---|

In [48]:
```python
geometry = [Point(xy) for xy in zip(fires_df['Longitude'], fires_df['Latitude'])]
geometry[:3]
crs = {'init': "EPSG:4326"}
fires_df1 = gpd.GeoDataFrame(fires_df, crs=crs, geometry=geometry)
fires_df1.head(2)
```

Out[48]:

| | Name | Latitude | Longitude | CaDate | CaYear | CaMonth | CaDay | ExitDate | AcresBurned |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Rim Fire | 37.857000 | -120.086000 | 2013-08-17 | 2013 | 8 | 17 | 2013-09-06 | 257314. |
| 1 | Powerhouse Fire | 34.585595 | -118.423176 | 2013-05-30 | 2013 | 5 | 30 | 2013-06-08 | 30274. |

Check for duplicate values based on geometry and date

In [49]:
```python
fires_df1 = fires_df1.sort_values(['CaDate'], ascending=True)
```

In [50]:
```python
fires_df1 = fires_df1[~fires_df1.duplicated(['geometry', 'CaDate'], keep='first')]
```

In [51]:
```python
fires_df1 = fires_df1[~fires_df1.duplicated(['Name','geometry', 'AcresBurned'], keep='first')]
```

In [52]:
```python
fires_df1 = fires_df1.sort_values('CaDate')
```

In [53]:
```python
fires_df1 = fires_df1.drop(["Archive_minus_Year"], axis = 1)
```

In [54]:
```python
fires_df1["ArchiveYear"].describe()
```

Out[54]:
```
count    1436.000000
mean     2016.804318
std         1.800893
min      2013.000000
25%      2016.000000
50%      2017.000000
75%      2018.000000
max      2019.000000
Name: ArchiveYear, dtype: float64
```

## Making Copies for the record

```
In [55]:   FireLocation = fire_location.copy()
           FirePolygon = ca_fires_df.copy()
           FireList = fires_df1.copy()
```

```
In [56]:   print(FirePolygon.crs)
           print(FireLocation.crs)
           print(FireList.crs)
```

```
epsg:4326
epsg:4326
+init=epsg:4326 +type=crs
```

```
In [57]:   FireLocation.crs
```

```
Out[57]:   <Geographic 2D CRS: EPSG:4326>
           Name: WGS 84
           Axis Info [ellipsoidal]:
           - Lat[north]: Geodetic latitude (degree)
           - Lon[east]: Geodetic longitude (degree)
           Area of Use:
           - name: World.
           - bounds: (-180.0, -90.0, 180.0, 90.0)
           Datum: World Geodetic System 1984 ensemble
           - Ellipsoid: WGS 84
           - Prime Meridian: Greenwich
```

```
In [58]:   FireList.crs
```

```
Out[58]:   <Geographic 2D CRS: +init=epsg:4326 +type=crs>
           Name: WGS 84
           Axis Info [ellipsoidal]:
           - lon[east]: Longitude (degree)
           - lat[north]: Latitude (degree)
           Area of Use:
           - name: World.
           - bounds: (-180.0, -90.0, 180.0, 90.0)
           Datum: World Geodetic System 1984 ensemble
           - Ellipsoid: WGS 84
           - Prime Meridian: Greenwich
```

# Merging Fire Perimeter with Fire Location by date.

Projecting the crs to from WGS84 to NAD83 so we can compute distances between points correctly in meters or kilometers.

```
In [59]:   FireLocation = FireLocation.to_crs({'init': "EPSG:3310"})
           FirePolygon = FirePolygon.to_crs({'init': "EPSG:3310"})
           FireList = FireList.to_crs({'init': "EPSG:3310"})
```

In [60]:
```python
FireLocation['Fire'] = 0
FireList['Fire'] = 0
```

In [61]:
```python
def get_nearestpoints(df1, df1day, df2, df2day, dist):
    """
        This Function merges dataframe for selected day by finding nearest points
    for each day and creates mini dfs for each day of month
    """
    days = list(range(1, 32))
    dfs = []
    for day in days:
        df = df1[df1[df1day] == day]
        df3 = df2
        m_df = gpd.sjoin_nearest(df, df3, how='left', max_distance = 10,distance_col=dist)
        m_df[dist] = m_df[dist].apply(lambda x: x/1000)
        d = pd.DataFrame(m_df)
        dfs.append(d)

    dfs = pd.concat(dfs)
    dfs = dfs.drop(["Fire"], axis = 1)
    return dfs
```

In [62]:
```python
def merge_data(data1, df1year, df1month, df1day, data2, df2year, df2month,df2day, dist, year):
    """
        This Function filters dataframe by year and months and calls for day dfs,
        append it and then converts it into pandas df.
    """
    months = list(range(1, 13))
    dfs = []
    for month in months:
        df1 = data1[(data1[df1year] == year) & (data1[df1month] == month)]
        df2 = data2[(data2[df2year] == year) & (data2[df2month] == month)]

        if 'Fire' in df2.columns:
            df = gpd.sjoin_nearest(df1, df2, how='left', distance_col=dist)
            df[dist] = df[dist].apply(lambda x: x/1000)
            d = pd.DataFrame(df)
            dfs.append(d)
        else:
            df = get_nearestpoints(df1, df1day, df2, df2day, dist)
            dfs.append(df)

    dfs = pd.concat(dfs)
    return dfs
```

In [63]:
```python
def get_data(df1, df1year, df1month, df1day, df2, df2year, df2month, df2day, d
ist):
    """
        This calls_ for all dataframes and combine it and create one dataset f
or
        fire data, so we can use the combined information to find the estimate
d dates
    """
    years = list(range(2011, 2021))
    dataframesList = []
    for year in years:
        data = merge_data(df1, df1year, df1month, df1day, df2, df2year, df2mon
th, df2day, dist, year)
        dataframesList.append(data)

    df = gpd.GeoDataFrame(pd.concat(dataframesList), crs=crs)
    try:
        df.drop('index_right', axis=1, inplace=True)
    except ValueError:
        # ignore if there are no index columns
        pass

    print(df.shape)

    return df
```

In [64]:
```python
nearestfire1 = get_data(FirePolygon, 'FireYear', 'FireMonth', 'FireDay',
                        FireLocation, 'DiscoveryYear', 'DiscoveryMonth', 'Disco
veryDay', 'locationdist')
```

```
(3838, 21)
```

# Handling Missing Values and Duplicates. Fixing Bad Data

### Analysis of Fire Start Dates Missing Values using Fire Location Data

In [65]:
```python
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
nearestfire1.head(2)
```

Out[65]:

| | ObjectID | State | UnitID | Name | FireCause | TotalAcres | geometry | FireDate | FireYear |
|---|---|---|---|---|---|---|---|---|---|
| **252** | 21694 | CA | SHF | FLAT | 10.0 | 62.040543 | MULTIPOLYGON (((-289859.162 314145.109, -28986... | 2020-06-30 | 2020.0 |
| **375** | 21819 | CA | SKU | NOYES 1-14 | 1.0 | 8.120131 | MULTIPOLYGON (((-224232.716 382534.303, -22423... | 2020-07-28 | 2020.0 |

2 rows × 21 columns

In [66]:
```python
def get_duration(df):
    df['LocDate_minus_FireDate'] = (df["DiscoveryDate"] - df["FireDate"]).dt.days
    df = df.sort_values('LocDate_minus_FireDate', ascending=True)
    return df
```

In [67]:
```python
nearestfire1 = get_duration(nearestfire1)
```

In [68]:
```python
# dropping duplicates based on earliest discovery date and location distance == 0 for fires bigger than 100 acres
nearestfire1 = nearestfire1[~((nearestfire1.index.duplicated(keep='first')) &
                             (nearestfire1['LocDate_minus_FireDate']<=0) &
                             (nearestfire1['TotalAcres']>100) & (nearestfire1['locationdist']==0))]

# dropping duplicates based on earliest discovery date and location distance == 0 for fires less than 100 acres
nearestfire1 = nearestfire1[~((nearestfire1.index.duplicated(keep='first')) &
                             (nearestfire1['LocDate_minus_FireDate']<=0) &
                             (nearestfire1['TotalAcres']<100) & (nearestfire1['locationdist']==0))]

# dropping it because it is not mapped correctly
nearestfire1 = nearestfire1[~((nearestfire1.index.duplicated(keep='first')) &
                             (nearestfire1['LocDate_minus_FireDate']>=0) &
                             (nearestfire1['TotalAcres']<100) & (nearestfire1['locationdist']==0))]

nearestfire1 = nearestfire1[~((nearestfire1.index.duplicated(keep='first')) &
                             (nearestfire1['LocDate_minus_FireDate']>=0) &
                             (nearestfire1['TotalAcres']>100) & (nearestfire1['locationdist']==0))]
```

In [69]:
```python
nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

```
In [70]: nearestfire1 = nearestfire1[~nearestfire1.index.duplicated(keep='first')]
         nearestfire1 = get_duration(nearestfire1)
         nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

```
In [71]: nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <0) &
                       (nearestfire1['DiscoveryDate'].notnull()) &
                       (nearestfire1['locationdist'] ==0)]

         ## replacing fire date with discovery date.
         nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <0) &
                          (nearestfire1['DiscoveryDate'].notnull()) &
                          (nearestfire1['locationdist'] ==0), 'FireDate'] = nearestfire
         1['DiscoveryDate']
```

```
In [72]: # replacing discovery dates for any fire which is less than kilometer away and
         starts earlier than reported fire date
         nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <0) &
                      (nearestfire1['DiscoveryDate'].notnull()) &
                      (nearestfire1['locationdist'] >0) &
                      (nearestfire1['locationdist'] <1) & (nearestfire1['TotalAcres']>1
         00)]

         nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <0) &
                          (nearestfire1['DiscoveryDate'].notnull()) &
                          (nearestfire1['locationdist'] >0) &
                          (nearestfire1['locationdist'] <1) & (nearestfire1['TotalAcre
         s']>100), 'FireDate'] = nearestfire1['DiscoveryDate']
```

```
In [73]: # replacing discovery dates for any fire which is less than kilometer away and
         starts earlier than reported fire date
         nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <0) &
                      (nearestfire1['DiscoveryDate'].notnull()) &
                      (nearestfire1['locationdist'] >0) &
                      (nearestfire1['locationdist'] <0.2) & (nearestfire1['TotalAcres']
         <100)]

         nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <0) &
                          (nearestfire1['DiscoveryDate'].notnull()) &
                          (nearestfire1['locationdist'] >0) &
                          (nearestfire1['locationdist'] <0.2) & (nearestfire1['TotalAcr
         es']<100), 'FireDate'] = nearestfire1['DiscoveryDate']
```

```
In [74]: nearestfire1 = get_duration(nearestfire1)
         nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

In [75]:
```python
# replacing discovery dates for any fire which is less than kilometer away and
starts earlier than reported fire date
nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
             (nearestfire1['DiscoveryDate'].notnull()) &
             (nearestfire1['LocDate_minus_FireDate'] >=-2) &
             (nearestfire1['locationdist'] <0.5)]

nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
             (nearestfire1['DiscoveryDate'].notnull()) &
             (nearestfire1['LocDate_minus_FireDate'] >=-2) &
             (nearestfire1['locationdist'] <0.5), 'FireDate'] = nearestfire1[
'DiscoveryDate']
```

In [76]:
```python
nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

In [77]:
```python
# replacing discovery dates for any fire which is less than kilometer away and
starts earlier than reported fire date
nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
                 (nearestfire1['DiscoveryDate'].notnull()) &
                 (nearestfire1['LocDate_minus_FireDate'] >=-2) &
                 (nearestfire1['locationdist'] <1),'FireDate'] = nearestfire1[
'DiscoveryDate']
```

In [78]:
```python
nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

In [79]:
```python
# replacing discovery dates for any fire which is less than 2 kilometer away a
nd
# starts earlier than reported fire date, but shares the same name
nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
             (nearestfire1['DiscoveryDate'].notnull()) &
             (nearestfire1['LocDate_minus_FireDate'] >=-2) &
             (nearestfire1['locationdist'] <2)]

nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
             (nearestfire1['DiscoveryDate'].notnull()) &
             (nearestfire1['LocDate_minus_FireDate'] >=-2) &
             (nearestfire1['locationdist'] <2),'FireDate'] = nearestfire1['Dis
coveryDate']

nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

In [80]:
```python
# replacing discovery dates for any fire which is less than 3 kilometer away a
nd
# starts earlier than reported fire date, but shares the same name
nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
            (nearestfire1['DiscoveryDate'].notnull()) &
            (nearestfire1['LocDate_minus_FireDate'] >=-3) &
            (nearestfire1['locationdist'] <=3)]

nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
            (nearestfire1['DiscoveryDate'].notnull()) &
            (nearestfire1['LocDate_minus_FireDate'] >=-3) &
            (nearestfire1['locationdist'] <3),'FireDate'] = nearestfire1['Dis
coveryDate']

nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

In [81]:
```python
# replacing discovery dates for any fire which is less than 10 kilometer away
 and
# starts earlier than reported fire date, but shares the same name and has hig
her totalacres
nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
            (nearestfire1['DiscoveryDate'].notnull()) &
            (nearestfire1['LocDate_minus_FireDate'] >=-10) &
            (nearestfire1['locationdist'] <=3) &
            (nearestfire1['TotalAcres']>=50)]

nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
            (nearestfire1['DiscoveryDate'].notnull()) &
            (nearestfire1['LocDate_minus_FireDate'] >=-10) &
            (nearestfire1['locationdist'] <=3) &
            (nearestfire1['TotalAcres']>=50),'FireDate'] = nearestfire1['Disc
overyDate']

nearestfire1 = get_duration(nearestfire1)
nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

```
In [82]:  # replacing discovery dates for any fire which is less than 10 kilometer away
          and
          # starts earlier than reported fire date, but shares the same name and has 100
          0+ totalacres
          nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
                       (nearestfire1['DiscoveryDate'].notnull()) &
                       (nearestfire1['LocDate_minus_FireDate'] >=-20) &
                       (nearestfire1['locationdist'] <=3) &
                       (nearestfire1['TotalAcres']>=1000)]

          nearestfire1.loc[(nearestfire1['LocDate_minus_FireDate'] <=-1) &
                       (nearestfire1['DiscoveryDate'].notnull()) &
                       (nearestfire1['LocDate_minus_FireDate'] >=-20) &
                       (nearestfire1['locationdist'] <=3) &
                       (nearestfire1['TotalAcres']>=1000),'FireDate'] = nearestfire1['Di
          scoveryDate']

          nearestfire1 = get_duration(nearestfire1)
          nearestfire1 = nearestfire1.sort_values('locationdist', ascending=True)
```

```
In [83]:  Fires_df1 = nearestfire1[['ObjectID', 'UnitID', 'FireCause', 'TotalAcres', 'ge
          ometry', 'FireDate', 'FireYear',
                                   'FireMonth', 'FireDay', 'Name']]

          Fires_df1[Fires_df1.duplicated(['geometry'], keep=False)]
```

Out[83]:

| ObjectID | UnitID | FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Na |
|----------|--------|-----------|------------|----------|----------|----------|-----------|---------|-----|
|  ◄ |  |  |  |  |  |  |  |  | ► |

```
In [84]:  Fires_df1.shape
```

Out[84]:  (3670, 10)

```
In [85]:  ## get left over fire location data using UniqueIdentifier to later merge into
          fire data.
          UniqueFireIdentifier1 = nearestfire1[(nearestfire1['LocDate_minus_FireDate'] <
          =-10) &
                                              (nearestfire1['DiscoveryDate'].notnull()) &
                                              (nearestfire1['locationdist'] >10)]
          UniqueFireIdentifier1 = UniqueFireIdentifier1[~UniqueFireIdentifier1.duplicate
          d(['UniqueFireIdentifier'], keep='first')]
          UniqueFireIdentifier1 = UniqueFireIdentifier1[['UniqueFireIdentifier']]
```

# Analysis of Fire Start Dates using State Fire List

```
In [86]:  nearestfire2 = get_data(Fires_df1, 'FireYear', 'FireMonth', 'FireDay',
                                 FireList, 'CaYear', 'CaMonth', 'CaDay', 'firedist')
```

          (3675, 24)

```
In [87]: def get_duration(df):
             df['CaDate_minus_FireDate'] = (df["CaDate"] - df["FireDate"]).dt.days
             df['Area_diff'] = (df["AcresBurned"] - df["TotalAcres"])
             df = df.sort_values('CaDate_minus_FireDate', ascending=True)

             return df
```

```
In [88]: nearestfire2 = get_duration(nearestfire2)
         nearestfire2 = nearestfire2.sort_values('firedist', ascending=True)
```

```
In [89]: # dropping duplicates based on earliest discovery date and location distance =
         = 0 for fires bigger than 100 acres
         nearestfire2 = nearestfire2[~((nearestfire2.index.duplicated(keep='first')) &
                                    (nearestfire2['CaDate_minus_FireDate']<=0) &
                                    (nearestfire2['TotalAcres']<100))]
```

```
In [90]: nearestfire2 = get_duration(nearestfire2)
         nearestfire2 = nearestfire2.sort_values('firedist', ascending=True)
```

```
In [91]: # dropping it because it is not mapped correctly
         nearestfire2 = nearestfire2[~((nearestfire2.index.duplicated(keep='first')) &
                                    (nearestfire2['CaDate_minus_FireDate']>=0) &
                                    (nearestfire2['TotalAcres']>100))]

         nearestfire2 = get_duration(nearestfire2)
         nearestfire2 = nearestfire2.sort_values('firedist', ascending=True)
```

```
In [92]: nearestfire2[nearestfire2.index.duplicated(keep=False)]
```

Out[92]:

| ObjectID | UnitID | FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Na |
|----------|--------|-----------|------------|----------|----------|----------|-----------|---------|-----|

0 rows × 26 columns

```
In [93]: ## No dates that are listed earlier for state fire list. We will use the origi
         nal fire date.
         nearestfire2[(nearestfire2['CaDate_minus_FireDate'] >-10) &
                      (nearestfire2['CaDate_minus_FireDate'] <=-1) &
                      (nearestfire2['CaDate'].notnull()) &
                      (nearestfire2['firedist'] <3)].sort_values('CaDate_minus_FireDat
         e', ascending=True)

         ## No dates that are listed earlier for state fire list. We will use the origi
         nal fire date.
         nearestfire2.loc[(nearestfire2['CaDate_minus_FireDate'] >-10) &
                          (nearestfire2['CaDate_minus_FireDate'] <=-1) &
                          (nearestfire2['CaDate'].notnull()) &
                          (nearestfire2['firedist'] <3), 'FireDate'] = nearestfire2['Ca
         Date']

         nearestfire2 = get_duration(nearestfire2)
         nearestfire2 = nearestfire2.sort_values('firedist', ascending=True)
```
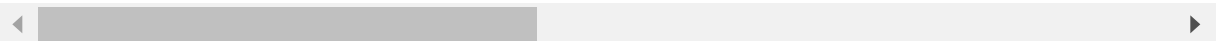
```
In [94]: UniqueId =  nearestfire2[(nearestfire2['CaDate_minus_FireDate'] >-20) &
                                 (nearestfire2['CaDate_minus_FireDate'] <=-1) &
                                 (nearestfire2['CaDate'].notnull()) &
                                 (nearestfire2['firedist'] >2)].sort_values('firedist'
         , ascending=True)
         UniqueId = UniqueId[~UniqueId.duplicated(['UniqueId'], keep='first')]
         UniqueId = UniqueId[['UniqueId']]
```

```
In [95]: nearestfire2.shape
```

```
Out[95]: (3670, 26)
```

```
In [96]: Fires_df2 = nearestfire2[['ObjectID', 'UnitID', 'FireCause', 'TotalAcres', 'ge
         ometry', 'FireDate', 'FireYear',
                                     'FireMonth', 'FireDay', 'Name_left']]

         Fires_df2[Fires_df2.duplicated(['geometry'], keep=False)]
```

Out[96]:

| ObjectID | UnitID | FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Na |
|---|---|---|---|---|---|---|---|---|---|
| ◄ | | | | | | | | | ► |

```
In [97]: UniqueFireIdentifier1 = UniqueFireIdentifier1['UniqueFireIdentifier'].tolist()
         UniqueId = UniqueId['UniqueId'].tolist()
```

```
In [98]: Fires_df4 = FireList[FireList['UniqueId'].isin(UniqueId)]
         Fires_df4 = Fires_df4[['geometry','CaDate', 'CaYear', 'CaMonth', 'CaDay', 'Acr
         esBurned', 'Name', 'UniqueId']]
```

```
In [99]: Fires_df4.shape
```

```
Out[99]: (334, 8)
```

## What is the Top Cause for Wildfire?

```
In [100]: Fires_df2['FireYear'].describe()
```

```
Out[100]: count    3670.000000
          mean     2015.915804
          std         2.875493
          min      2011.000000
          25%      2013.000000
          50%      2016.000000
          75%      2018.000000
          max      2020.000000
          Name: FireYear, dtype: float64
```

```
In [101]: Fires_df2 = Fires_df2[Fires_df2['FireYear'] <=2020]
```

In [102]:
```python
Fires_df2["FireCause"].replace({1: "Lightning", 2: "Equipment Use", 3:"Smokin
g",
                                4:"Campfire", 5:"Debris", 6:"Railroad", 7:"Ars
on", 8:"Playing with Fire",
                                9:"Miscellaneous",10:"Vehicle", 11:"Powerline"
, 14:"Unknown/ Unidentified",
                                15: 'Structure', 16: 'Aircraft', 18: 'Escaped
 Prescribed Burn', 19: "Illegal Alien Campfire"}, inplace=True)
```

In [103]:
```python
plt.rcParams['figure.figsize'] = [10,6]
sns.set(font_scale = 1.3)
sns.set_style("white")
df2 = Fires_df2[Fires_df2['TotalAcres'] >99]
ax = sns.countplot(y=df2['FireCause'],data=df2, palette="flare",
                   order=df2['FireCause'].value_counts().index)

ax.set_title("Cause of Big Wildfires from 2011-2020",fontsize = 15, loc='left'
)
ax.set_xlabel("Count")
ax.set_ylabel("Fire Cause")
sns.despine()

plt.show()
```



In [104]:
```python
Fires_df2 = Fires_df2.drop(['ObjectID', 'UnitID'], axis = 1)
```

In [105]:
```python
Fires_df2.rename(columns={"Name_left":"Name"}, inplace=True)
Fires_df2['UniqueId'] = "0B"
```

In [106]:
```python
Fires_df4.rename(columns={"CaDate":"FireDate", "CaYear":"FireYear",
                          "CaMonth":"FireMonth", "CaDay":"FireDay", "AcresBurn
ed":"TotalAcres"}, inplace=True)
```

In [107]: `Fires_df4[Fires_df4.duplicated(['geometry'], keep=False)]`

Out[107]:

| geometry | FireDate | FireYear | FireMonth | FireDay | TotalAcres | Name | UniqueId |
|---|---|---|---|---|---|---|---|

In [108]: `Fires_df = Fires_df2.append(Fires_df4)`

In [109]: `Fires_df = Fires_df[~(Fires_df['TotalAcres'] == 0)]`

In [110]: `Fires_df = Fires_df.sort_index()`

In [111]: `Fires_df['FireYear'].describe()`

Out[111]:
```
count    3995.000000
mean     2015.954693
std         2.810346
min      2011.000000
25%      2013.500000
50%      2017.000000
75%      2018.000000
max      2020.000000
Name: FireYear, dtype: float64
```
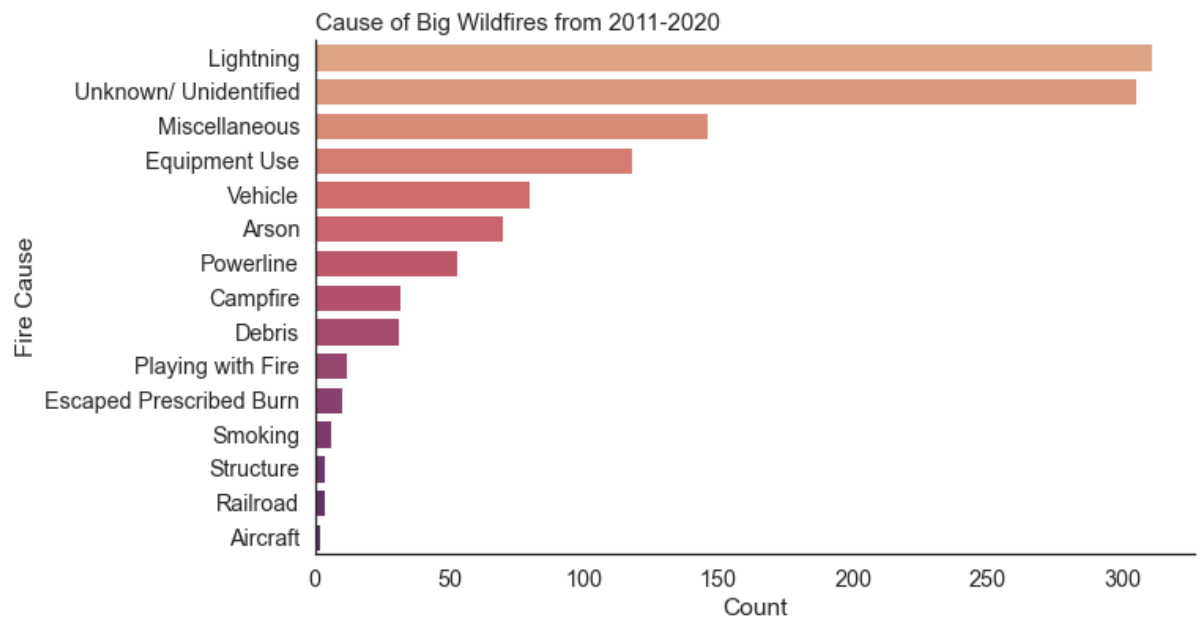
In [112]: `Fires_df = Fires_df`

In [113]: `Fires_df[['TotalAcres']].describe()`

Out[113]:

| | TotalAcres |
|---|---|
| count | 3.988000e+03 |
| mean | 2.983130e+03 |
| std | 2.536381e+04 |
| min | 1.356887e-03 |
| 25% | 1.060945e+01 |
| 50% | 3.700000e+01 |
| 75% | 2.032367e+02 |
| max | 1.032699e+06 |

**Check for duplicates by Name**

In [114]:
```python
import string
import re

def text_w_punc(text):
    pattern = r'[^A-Za-z ]'
    if re.search("[^0-9]", text):
        pass
    else:
        regex = re.compile(pattern)
        text = regex.sub(' ', text)
    return text

def no_whitespace(string):
    string = string.replace('fire','')
    return string.strip()

def string_lower(df, col):
    # converting all text to lowercase
    df[col] = df[col].str.lower()
    return df
```

In [115]:
```python
Fires_df['Name'] = Fires_df['Name'].astype(str)
Fires_df['Name'] = Fires_df['Name'].apply(no_whitespace)
Fires_df['Name'] = Fires_df['Name'].apply(lambda x: "".join([i for i in x if i
not in string.punctuation]))
Fires_df['Name'] = Fires_df['Name'].apply(lambda x: text_w_punc(x))
Fires_df = string_lower(Fires_df, "Name")
Fires_df['Name'] = Fires_df['Name'].apply(no_whitespace)
```

In [116]:
```python
Fires_df.shape
```

Out[116]:
```
(3995, 9)
```

In [117]:
```python
f = Fires_df[(Fires_df.duplicated(['Name', 'FireDate'], keep=False)) &
             (Fires_df['Name']!="")].sort_values(['Name','FireDate', 'TotalAcr
es'], ascending=False)
print(f.shape)

Fires_df = Fires_df[~((Fires_df.duplicated(['Name', 'FireDate'], keep=False))
&
                      (Fires_df['Name']!=""))]
print(Fires_df.shape)
```

```
(415, 9)
(3580, 9)
```

In [118]:
```python
multipolygon =f[f['UniqueId']=="0B"]
f = f[~(f['UniqueId']=="0B")]
```

In [119]:
```python
Fires_df = Fires_df.append(multipolygon)
```

```
In [120]:  Fires_df.shape
```

Out[120]:  (3803, 9)

```
In [121]:  Fires_df[Fires_df.duplicated(['geometry'], keep=False)]
```

Out[121]:

| FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Name | UniqueId |
|-----------|------------|----------|----------|----------|-----------|---------|------|----------|

## What fire class is more common?

```
In [122]:  ## Fire Class Binning
           # A=greater than 0 but less than or equal to 0.25 acres
           # B=0.26-9.9 acres, C=10.0-99.9 acres, D=100-299 acres
           # E=300 to 999 acres, F=1000 to 4999 acres, and G=5000+ acres)

           #binning method for confidence of fire.
           bins = [0,.25,9.9,99.9,299,999,4999,1032699.0]
           labels = ['A', 'B','C', 'D', 'E', 'F', 'G']
           Fires_df['FireSize'] = pd.cut(Fires_df['TotalAcres'], bins=bins, labels=labels
           )
           Fires_df['FireSize']= Fires_df['FireSize'].fillna('A')
```

In [123]:
```python
## Fire Class Binning
# A=greater than 0 but less than or equal to 0.25 acres
# B=0.26-9.9 acres, C=10.0-99.9 acres, D=100-299 acres
# E=300 to 999 acres, F=1000 to 4999 acres, and G=5000+ acres)

plt.rcParams['figure.figsize'] = [10,6]
colors = ["#FF800D", "#E8540C","#FF3000", "#C5140C", "#7A3000", "#6E1818", "#3
B0D0D"]
sns.set_palette(sns.color_palette(colors))
sns.set(font_scale = 1.3)
sns.set_style("white")


ax = sns.countplot(x=Fires_df['FireSize'],data=Fires_df, palette=sns.color_pal
ette(colors),
                    order=Fires_df['FireSize'].value_counts().index)

ax.set_title("Number of Wildfires by Sizes from 2011-2020",fontsize = 15, loc=
'left')
ax.set_xlabel("Fire Class")
ax.set_ylabel("Count")
sns.despine()

plt.show()
```

Number of Wildfires by Sizes from 2011-2020

```python
In [124]: import textwrap
          def wrap_labels(ax, width, break_long_words=False):
              labels = []
              for label in ax.get_xticklabels():
                  text = label.get_text()
                  labels.append(textwrap.fill(text, width=width,
                                break_long_words=break_long_words))
              ax.set_xticklabels(labels, rotation=0)

          def plot_cause(df, title):


              plt.rcParams['figure.figsize'] = [10,6]
              sns.set(font_scale = 1.3)
              sns.set_style("white")

              ax = sns.countplot(x=df['FireCause'],data=df, palette="flare",
                             order=df['FireCause'].value_counts().iloc[:5].index)
              ax.set_title(title,fontsize = 15, loc='left')
              ax.set_xticklabels(ax.get_xticklabels(), rotation=90)
              #ax.tick_params(axis='x', labelrotation=90)
              wrap_labels(ax, 5)
              ax.figure
              ax.set_xlabel("Fire Cause")
              ax.set_ylabel("Count")
              sns.despine()

              plt.show()
```

## What are the Top Causes for Small Wildfires (less than 300 Acres)?

```
In [125]: df = Fires_df[(Fires_df['TotalAcres'] <300)]
          plot_cause(df, "Top Causes for Fires Less than 300 Acres")
```

Top Causes for Fires Less than 300 Acres



```
In [126]: print("Number of Small fires: {}".format(df.shape[0]))
```

Number of Small fires: 3017

## What are the Top Causes for Large Wildfires (greater than or 300 Acres)?

In [127]:
```python
df = Fires_df[(Fires_df['TotalAcres'] >=300)]
plot_cause(df, "Top Causes for Fires More than 300 Acres")
```

Top Causes for Fires More than 300 Acres



In [128]:
```python
print("Number of Large fires: {}".format(df.shape[0]))
```

Number of Large fires: 779

In [129]:
```python
Datatype(Fires_df)
```

```
There are 3803 rows and 10 columns
FireCause              object
TotalAcres            float64
geometry             geometry
FireDate        datetime64[ns]
FireYear              float64
FireMonth             float64
FireDay               float64
Name                   object
UniqueId               object
FireSize             category
dtype: object
['TotalAcres' 'FireYear' 'FireMonth' 'FireDay']
['FireCause' 'geometry' 'FireDate' 'Name' 'UniqueId' 'FireSize']
```

In [130]:
```python
geo_fires_df = Fires_df
geo_fires_df[geo_fires_df.duplicated(['geometry'], keep=False)]
```

Out[130]:

| FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Name | UniqueId | Fir |
|-----------|------------|----------|----------|----------|-----------|---------|------|----------|-----|

## 2017 had the highest number of wildfire

```
In [131]: #ca_fires_df = ca_fires_df[ca_fires_df["ReportedAcres"] >= .25]

          fire_count = pd.DataFrame(geo_fires_df['FireYear'].value_counts(sort=False))

          plt.rcParams['figure.figsize'] = [10,6]
          ax = sns.lineplot(data=fire_count, x=fire_count.index, y="FireYear", color="#F
          F3000", linewidth = 2)


          ax.set_title("Total Fires in California (2011-2020)",fontsize = 15,loc='left')
          ax.set_xlabel("Year")
          ax.set_ylabel("Count")
          sns.despine()
          plt.show()
```



## Which Month had the highest fire counts over time?

```
In [132]: sns.set(font_scale = 1.3)
          sns.set_style("white")
          df1 = geo_fires_df.replace({'FireMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar', 4:
          'Apr',
                                                      5: 'May', 6: 'Jun', 7: 'Jul', 8: 'Au
          g',
                                                      9: 'Sep', 10: 'Oct', 11: 'Nov', 12:
          'Dec'}})
          ax = sns.countplot(df1['FireMonth'],data=df1, color="#FF3000",
                             order=df1['FireMonth'].value_counts().index)

          ax.set_title("Total Fires for Each Month (2011-2020)",fontsize = 15, loc='lef
          t')
          ax.set_xlabel("Month")
          ax.set_ylabel("Count")
          sns.despine()

          plt.show()
```



Total Fires for Each Month (2011-2020)

```
In [133]: plot_gf1 = geo_fires_df.to_crs({'init': "EPSG:4326"})
```

## What area of California is more wildfire prone area?

In [134]:
```python
fig, ax = plt.subplots(figsize = (10,10))
fig.suptitle('Mapped California Fire Data (2010-2020)', fontsize=15)
plt.xticks([-124, -122, -120, -118, -116, -114])

USA[USA.STATEFP == '06'].plot(ax = ax, edgecolor="darkgrey", facecolor='lightg
rey')
plot_gf1.plot(ax=ax, color="#FF3000")
```

Out[134]: <AxesSubplot:>



Mapped California Fire Data (2010-2020)

In [135]: `geo_fires_df.crs`

Out[135]: 
```
<Derived Projected CRS: EPSG:3310>
Name: NAD83 / California Albers
Axis Info [cartesian]:
- E[east]: Easting (metre)
- N[north]: Northing (metre)
Area of Use:
- name: United States (USA) - California.
- bounds: (-124.45, 32.53, -114.12, 42.01)
Coordinate Operation:
- name: California Albers
- method: Albers Equal Area
Datum: North American Datum 1983
- Ellipsoid: GRS 1980
- Prime Meridian: Greenwich
```

In [136]:
```python
def add(date, num):
    end_date = date + dt.timedelta(days=num)
    return end_date

def substract(date, num):
    end_date = date - dt.timedelta(days=num)
    return end_date
```

In [137]:
```python
Dates = geo_fires_df[['FireDate']]
Dates = Dates.drop_duplicates()
```

In [138]:
```python
Dates = Dates['FireDate'].tolist()
```

In [139]:
```python
years = list(range(2010, 2021))
months = list(range(1, 13))
days = list(range(1, 32))
```

In [140]:
```python
fires_copy = geo_fires_df
```

In [142]:
```python
fires_copy.shape
```

Out[142]: `(3803, 10)`

In [ ]:

# Appendix A.2 MODIS Data Preparation code

In [2]: ▶|

```python
1   import datetime as dt
2   from pathlib import Path
3   import math
4   import os
5   import sqlite3
6   import json
7   import geopandas as gpd
8   import pygeos
9   import pyproj
10  import shapely
11  import shapely.ops as ops
12  from shapely.geometry import Point, Polygon
13  from shapely.geometry.polygon import Polygon
14  from functools import partial
15
16  import pandas as pd
17  import numpy as np
18  import seaborn as sns
19  import matplotlib.pyplot as plt
20  %matplotlib inline
21
22  from sklearn.model_selection import train_test_split
23
24
25  from sklearn import svm
26  from sklearn.svm import SVC
27  from sklearn.ensemble import RandomForestClassifier
28  from sklearn.naive_bayes import GaussianNB
29  from sklearn.metrics import accuracy_score, classification_report, confu
30
31  from sklearn.feature_selection import SelectKBest
32  from sklearn.feature_selection import chi2, f_classif, mutual_info_class
33  from functools import partial
34
35
36  from sklearn.preprocessing import StandardScaler
37
38  import warnings
39  warnings.filterwarnings('ignore')
```

## Data Collection

In [3]: ▶|

```python
1   ## Map file
```

In [4]:
```python
1  USA = gpd.read_file("Data/County/cb_2018_us_county_500k.shp")
2  USA.head()
```

Out[4]:

| | STATEFP | COUNTYFP | COUNTYNS | AFFGEOID | GEOID | NAME | LSAD | ALAND |
|---|---|---|---|---|---|---|---|---|
| **0** | 21 | 007 | 00516850 | 0500000US21007 | 21007 | Ballard | 06 | 639387454 |
| **1** | 21 | 017 | 00516855 | 0500000US21017 | 21017 | Bourbon | 06 | 750439351 |
| **2** | 21 | 031 | 00516862 | 0500000US21031 | 21031 | Butler | 06 | 1103571974 |
| **3** | 21 | 065 | 00516879 | 0500000US21065 | 21065 | Estill | 06 | 655509930 |
| **4** | 21 | 069 | 00516881 | 0500000US21069 | 21069 | Fleming | 06 | 902727151 |

**2b. Load the dataset 2: NASA Active Fire Data - https://earthdata.nasa.gov/ (https://earthdata.nasa.gov/)**

In [5]:
```python
1  current_dir = Path(os.getcwd()).absolute()
```

```
In [6]:    1  old_nasa_df = pd.read_csv('Data/fire_archive_M6_156000.csv') # archive da
           2  new_nasa_df = pd.read_csv('Data/fire_nrt_M6_156000.csv') # new data
           3
           4  nasa_df = pd.concat([old_nasa_df, new_nasa_df]) #concatenate old and new
           5  print(nasa_df.shape)
           6  nasa_df.tail(4)
```

(1248606, 15)

Out[6]:

| | latitude | longitude | brightness | scan | track | acq_date | acq_time | satellite | instrument |
|---|---|---|---|---|---|---|---|---|---|
| **66330** | 36.293 | -118.664 | 329.6 | 1.1 | 1.0 | 2020-09-27 | 1855 | Terra | MODIS |
| **66331** | 40.043 | -123.080 | 337.5 | 1.2 | 1.1 | 2020-09-27 | 1855 | Terra | MODIS |
| **66332** | 45.726 | -118.308 | 305.5 | 1.0 | 1.0 | 2020-09-27 | 1855 | Terra | MODIS |
| **66333** | 37.262 | -119.443 | 319.4 | 1.0 | 1.0 | 2020-09-27 | 1855 | Terra | MODIS |

## Data Preliminary Analysis

```
In [8]:    1  # check for missing value
           2  def percentMissing(df):
           3
           4      df_numeric = df.select_dtypes(include=[np.number])
           5      numeric_cols = df_numeric.columns.values
           6
           7      # % of missing data
           8      for col in df.columns:
           9          # create missing indicator for features with missing data
          10          missing = df[col].isnull()
          11          pct_missing = np.mean(missing)*100
          12          #if pct_missing >60:
          13          print('{} - {}%'.format(col, round(pct_missing)))
          14          num_missing = np.sum(missing)
```

In [9]: ▶|

```python
1  # Checking data type
2  def Datatype(df):
3      # shape and data types of the data
4      print("There are {} rows and {} columns".format(df.shape[0], df.shape[1]))
5      print(df.dtypes)
6
7      # select numeric columns
8      df_numeric = df.select_dtypes(include=[np.number])
9      numeric_cols = df_numeric.columns.values
10     print(numeric_cols)
11
12     # select non numeric columns
13     df_non_numeric = df.select_dtypes(exclude=[np.number])
14     non_numeric_cols = df_non_numeric.columns.values
15     print(non_numeric_cols)
```

# Data Exploration: MODIS Collection 6 Active Fire Data

In [10]: ▶|

```python
1  Datatype(nasa_df)
```

```
There are 1248606 rows and 15 columns
latitude       float64
longitude      float64
brightness     float64
scan           float64
track          float64
acq_date        object
acq_time         int64
satellite       object
instrument      object
confidence       int64
version         object
bright_t31     float64
frp            float64
daynight        object
type           float64
dtype: object
['latitude' 'longitude' 'brightness' 'scan' 'track' 'acq_time'
 'confidence' 'bright_t31' 'frp' 'type']
['acq_date' 'satellite' 'instrument' 'version' 'daynight']
```

In [11]: ▶|

```python
1  # Adding new month and day variables
2  nasa_df['acq_date'] = pd.to_datetime(nasa_df['acq_date'])
3  nasa_df.rename(columns={"acq_date":"ActiveDate"}, inplace=True)
4
5  nasa_df['ActiveYear'] = nasa_df['ActiveDate'].dt.year
6  nasa_df['ActiveMonth'] = nasa_df['ActiveDate'].dt.month
7  nasa_df['ActiveDay'] = nasa_df['ActiveDate'].dt.day
```

In [12]:
```python
#binning method for confidence of fire.
bins = [0, 30,80,100]
labels = ['low', 'nominal','high']
nasa_df['ConfidenceBinned'] = pd.cut(nasa_df['confidence'], bins=bins, l
nasa_df['ConfidenceBinned']= nasa_df['ConfidenceBinned'].fillna('low')
```

In [13]:
```python
# dropping version and instrument variable because it just tells us what
nasa_df = nasa_df.drop(['instrument', 'version', 'acq_time'], axis = 1)
nasa_df = nasa_df.rename(columns={'brightness': 'Brightness', 'scan': 'S
                                  'track': 'Track', 'longitude': 'NasaLor
                                  'satellite': 'Satellite', 'confidence'
                                  'bright_t31':'BrightT31', 'frp': 'Frp'
                                  'type': 'HotSpotType','latitude': 'Nasa

nasa_df.shape
```

Out[13]: (1248606, 16)

In [14]:
```python
nasa_df.dropna(inplace=True)
```

In [15]:
```python
ca_nasa_df = nasa_df[(nasa_df['NasaLatitude']<= 42) & (nasa_df['NasaLati
ca_nasa_df = ca_nasa_df[(ca_nasa_df['NasaLongitude']<= -114) & (ca_nasa_
```

In [16]:
```python
x = ca_nasa_df[ca_nasa_df['NasaLatitude']<= 42]
y = ca_nasa_df[(ca_nasa_df['NasaLatitude'] >= 42) & (ca_nasa_df['NasaLat
y = y[y['NasaLongitude'] <=-117.1]
```

In [17]:
```python
ca_nasa_df = ca_nasa_df[ca_nasa_df['ActiveYear'] >=2011]
```

In [18]:
```python
ca_nasa_df.shape
```

Out[18]: (114599, 16)

In [19]:
```python
percentMissing(ca_nasa_df)
```

```
NasaLatitude - 0%
NasaLongitude - 0%
Brightness - 0%
Scan - 0%
Track - 0%
ActiveDate - 0%
Satellite - 0%
Confidence - 0%
BrightT31 - 0%
Frp - 0%
DayNight - 0%
HotSpotType - 0%
ActiveYear - 0%
ActiveMonth - 0%
ActiveDay - 0%
ConfidenceBinned - 0%
```

In [20]: ▶

```
1  # check for duplicates in coordinates
2  duplicate = ca_nasa_df[ca_nasa_df.duplicated(['NasaLatitude', 'NasaLongi
3  duplicate
```

Out[20]:

| | NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Satellite | Confiden |
|---|---|---|---|---|---|---|---|---|
| **587108** | 36.8276 | -118.8827 | 309.6 | 2.6 | 1.5 | 2015-08-20 | Terra | |
| **589099** | 36.8276 | -118.8827 | 406.3 | 1.7 | 1.3 | 2015-08-20 | Aqua | 1 |
| **701970** | 36.3855 | -121.7702 | 306.1 | 3.5 | 1.8 | 2016-07-30 | Terra | |
| **702125** | 36.3855 | -121.7702 | 373.7 | 1.0 | 1.0 | 2016-07-30 | Terra | 1 |
| **982192** | 40.7501 | -122.5168 | 328.5 | 1.1 | 1.0 | 2018-08-06 | Terra | 1 |
| **982434** | 40.7501 | -122.5168 | 321.9 | 1.0 | 1.0 | 2018-08-06 | Aqua | 1 |

**They are not really duplicates and got detected by two separate Satellite**

In [21]: ▶

```
1  ca_nasa_df.shape
```

Out[21]: (114599, 16)

In [22]: ▶

```
1  # Histograms
2  def histogram(xaxes, yaxes, df, x, y, nrows, color):
3      plt.rcParams['figure.figsize'] = (x, y)
4
5      fig, axes = plt.subplots(nrows = nrows, ncols = 2)
6      fig.suptitle('Distribution of Fire Pixel Attributes in West Coast Reg
7
8      # draw histograms in for loop
9      axes = axes.ravel()
10     for idx, ax in enumerate(axes):
11         # drops NaN values
12         ax.hist(df[num_features[idx]].dropna(), bins=40, color= color)
13         ax.set_xlabel(xaxes[idx], fontsize=15)
14         ax.set_ylabel(yaxes[idx], fontsize=15)
15         ax.tick_params(axis='both', labelsize=15)
16         right_side = ax.spines["right"]
17         right_side.set_visible(False)
18         top = ax.spines["top"]
19         top.set_visible(False)
20
21     plt.show()
```

In [23]: ▶|
```python
1  # Specify the features of interest
2  num_features = ['Brightness', 'Scan', 'Frp', 'BrightT31']
3  xaxes = num_features
4  yaxes = ['Counts', 'Counts', 'Counts', 'Counts']
5  histogram(xaxes, yaxes, ca_nasa_df, 15,8, 2, "#ffa800")
```

Distribution of Fire Pixel Attributes in West Coast Region (2011-2020)



In [24]: ▶|
```python
1  geometry = [Point(xy) for xy in zip(ca_nasa_df['NasaLongitude'], ca_nasa_
2  geometry[:3]
```

Out[24]: [<shapely.geometry.point.Point at 0x1f833064220>,
          <shapely.geometry.point.Point at 0x1f82365e740>,
          <shapely.geometry.point.Point at 0x1f81d7ff370>]

In [25]: ▶|
```python
1  crs = {'init': "EPSG:4326"}
2  geo_nasa_df = gpd.GeoDataFrame(ca_nasa_df, crs=crs, geometry=geometry)
3  geo_nasa_df.head(2)
```

Out[25]:

| | NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Satellite | Confidenc |
|---|---|---|---|---|---|---|---|---|
| **26628** | 34.6033 | -118.3848 | 305.8 | 1.2 | 1.1 | 2011-01-06 | Terra | 5 |
| **26638** | 39.8467 | -121.5210 | 308.2 | 1.1 | 1.0 | 2011-01-07 | Terra | 7 |

In [26]: ▶|
```python
1  geo_nasa_df = geo_nasa_df.to_crs({'init': "EPSG:3310"})
```

In [27]:

```python
y = ca_nasa_df[ca_nasa_df['NasaLongitude'] <=-119.5]
z = ca_nasa_df[(ca_nasa_df['NasaLatitude'] <= 39.5) & (ca_nasa_df['NasaL
w = z[(z['NasaLongitude'] >=-119.5) & (z['NasaLongitude'] <=-116)]
v = ca_nasa_df[ca_nasa_df['NasaLatitude'] <= 36.5]

plot_df = pd.concat([y,z,w,v])
plot_df = plot_df[~((plot_df['NasaLatitude']>39) & (plot_df['NasaLongitu
plot_df = plot_df[~((plot_df['NasaLatitude']>38) & (plot_df['NasaLongitu
plot_df = plot_df[~((plot_df['NasaLatitude']>37) & (plot_df['NasaLongitu
plot_df = plot_df[~((plot_df['NasaLatitude']>36) & (plot_df['NasaLongitu


geometry = [Point(xy) for xy in zip(plot_df['NasaLongitude'], plot_df['Na
geometry[:3]
plot_df = gpd.GeoDataFrame(plot_df, crs=crs, geometry=geometry)
```

In [28]:

```python
stcode = ['06']
```

In [29]: ▶|
```
1  fig, ax = plt.subplots(figsize = (10, 10))
2  fig.suptitle('Geospatial Plot Calfornia Fire Pixels (2011-2020)', fontsi
3
4  plot_df[plot_df['ConfidenceBinned'] == "low"].plot(ax=ax, color="#FFDF0D"
5  plot_df[plot_df['ConfidenceBinned'] == "nominal"].plot(ax=ax, color="#FF/
6  plot_df[plot_df['ConfidenceBinned'] == "high"].plot(ax=ax, color="#E85400
7  USA[USA['STATEFP'].isin(stcode)].plot(ax=ax, edgecolor="grey", facecolor
8  plt.legend(prop={'size':15})
```

Out[29]:  <matplotlib.legend.Legend at 0x1f832f50820>



Geospatial Plot Calfornia Fire Pixels (2011-2020)

In [30]: ⏭

```
 1  confidence_count = pd.DataFrame(geo_nasa_df[['ConfidenceBinned', 'Active'
 2  confidence_count.columns.values[2] = 'count'
 3  colors = ["#FFDF0D", "#FFA800","#E8540C"]
 4  plt.rcParams['figure.figsize'] = [10,6]
 5  ax = sns.lineplot(data=confidence_count, x="ActiveYear", y='count', hue=
 6                    palette=sns.color_palette(colors), linewidth=2)
 7
 8  ax.set_title("Total Nominal-High Confidence Fire Detected in California A
 9  ax.set_xlabel("Year")
10  ax.set_ylabel("Confidence fire")
11  sns.despine()
12  plt.legend(prop={'size':15}, loc='upper left')
13  plt.show()
14
```



## What Hotspot Type has the highest confidence for Fire?

In [31]:

```
1  colors = ["#FFDF0D", "#FFA800","#FF710D"]
2  df1 = geo_nasa_df.replace({'HotSpotType' : {3 : 'Offshore', 2 : 'Other S
3                                              1 : "Active volcano", 0 : "Presu
4  sns.set_palette(sns.color_palette(colors))
5  plt.rcParams['figure.figsize'] = [10,6]
6  sns.countplot(y=df1['HotSpotType'], data=df1, hue=df1['ConfidenceBinned'
7                                      xlabel = "Count", ylabel = "Area Ty
8
9  sns.despine()
10 plt.legend(prop={'size':15}, loc='lower right')
11 plt.show()
```



## What Month highest fires were detected?

In [32]:

```python
colors = ["#FFDF0D", "#FFA800","#FF710D"]
df1 = geo_nasa_df.replace({'ActiveMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar
                                            5: 'May', 6: 'Jun', 7: 'Jul', 8
                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
sns.set_palette(sns.color_palette(colors))
plt.rcParams['figure.figsize'] = [10,6]
sns.countplot(x=df1['ActiveMonth'], data=df1, hue=df1['ConfidenceBinned']
              order=df1['ActiveMonth'].value_counts().index).set(title =
                                            xlabel = "Month", ylabel = "Count")

sns.despine()
plt.legend(prop={'size':15}, loc='upper right')
plt.show()
```



Confidence of Fire Pixel by Month

# Appendix A3: Soil and Metereological Data Preparation

In [1]:
```python
1   import datetime as dt
2   from pathlib import Path
3   import math
4   import os
5   import sqlite3
6   import json
7   import geopandas as gpd
8   import pygeos
9   import pyproj
10  import shapely
11  import shapely.ops as ops
12  from shapely.geometry import Point, Polygon
13  from shapely.geometry.polygon import Polygon
14  from functools import partial
15
16  import pandas as pd
17  import numpy as np
18  import seaborn as sns
19  import matplotlib.pyplot as plt
20  %matplotlib inline
21
22  from sklearn.model_selection import train_test_split
23
24
25  from sklearn import svm
26  from sklearn.svm import SVC
27  from sklearn.ensemble import RandomForestClassifier
28  from sklearn.naive_bayes import GaussianNB
29  from sklearn.metrics import accuracy_score, classification_report, confu
30
31  from sklearn.feature_selection import SelectKBest
32  from sklearn.feature_selection import chi2, f_classif, mutual_info_class
33  from functools import partial
34
35
36  from sklearn.preprocessing import StandardScaler
37
38  import warnings
39  warnings.filterwarnings('ignore')
```

## Data Collection

In [2]:
```python
1   Drought1 = pd.read_csv('Data/DroughtData2010-2011.csv') # archive data
2   Drought2 = pd.read_csv('Data/DroughtData2012-2020.csv') # new data Drougl
3   Drought3 = pd.read_csv('Data/DroughtData2000.csv')
```

```
In [3]:    1  soil_df = pd.read_csv('Data/soil_data.csv') # new data
```

```
In [4]:    1  # Checking data type
           2  def Datatype(df):
           3      # shape and data types of the data
           4      print("There are {} rows and {} columns".format(df.shape[0], df.shape
           5      print(df.dtypes)
           6
           7      # select numeric columns
           8      df_numeric = df.select_dtypes(include=[np.number])
           9      numeric_cols = df_numeric.columns.values
          10      print(numeric_cols)
          11
          12      # select non numeric columns
          13      df_non_numeric = df.select_dtypes(exclude=[np.number])
          14      non_numeric_cols = df_non_numeric.columns.values
          15      print(non_numeric_cols)
```

## 1c. Data Preliminary Analysis: Drought Data (2010-2020)

```
In [5]:    1  # Adding new month and day variables
           2  Drought3['date'] = pd.to_datetime(Drought3['date'])
           3  Drought3.rename(columns={"date":"DroughtDate"}, inplace=True)
           4
           5  Drought3['DroughtYear'] = Drought3['DroughtDate'].dt.year
           6  Drought3['DroughtMonth'] = Drought3['DroughtDate'].dt.month
           7  Drought3['DroughtDay'] = Drought3['DroughtDate'].dt.day
```

```
In [6]:    1  df1 = Drought3[Drought3['DroughtYear'] >=2011]
```

```
In [7]:    1  df1['DroughtYear'] = df1['DroughtYear'].astype(int)
           2  df1['DroughtMonth'] = df1['DroughtMonth'].astype(int)
           3  df1['DroughtDay'] = df1['DroughtDay'].astype(int)
```

```
In [8]:  1  drought_df = pd.concat([Drought1, Drought2]) #concatenate old and new da
         2  print(drought_df.shape)
         3  drought_df.tail(4)
```

(4540788, 21)

Out[8]:

|  | fips | date | PRECTOT | PS | QV2M | T2M | T2MDEW | T2MWET | T2M_MAX | T2M_M |
|---|---|---|---|---|---|---|---|---|---|---|
| **2271944** | 56043 | 2020-12-28 | 0.00 | 83.04 | 1.82 | -7.31 | -12.06 | -9.68 | -1.48 | -11. |
| **2271945** | 56043 | 2020-12-29 | 0.00 | 82.78 | 1.87 | -7.38 | -11.79 | -9.59 | -0.88 | -11. |
| **2271946** | 56043 | 2020-12-30 | 0.01 | 82.87 | 1.57 | -6.40 | -13.94 | -10.17 | 1.33 | -12. |
| **2271947** | 56043 | 2020-12-31 | 0.00 | 82.82 | 2.13 | -3.83 | -10.12 | -6.98 | 2.16 | -8. |

4 rows × 21 columns

```
In [9]:  1  # Adding new month and day variables
         2  drought_df['date'] = pd.to_datetime(drought_df['date'])
         3  drought_df.rename(columns={"date":"DroughtDate"}, inplace=True)
         4
         5  drought_df['DroughtYear'] = drought_df['DroughtDate'].dt.year
         6  drought_df['DroughtMonth'] = drought_df['DroughtDate'].dt.month
         7  drought_df['DroughtDay'] = drought_df['DroughtDate'].dt.day
```

```
In [10]:  1  drought_df['DroughtYear'] = drought_df['DroughtYear'].astype(int)
          2  drought_df['DroughtMonth'] = drought_df['DroughtMonth'].astype(int)
          3  drought_df['DroughtDay'] = drought_df['DroughtDay'].astype(int)
```

```
In [11]:  1  df = pd.concat([drought_df, df1]) #concatenate old and new data
          2  print(df.shape)
          3  df.tail(4)
```

(11353524, 24)

Out[11]:

|  | fips | DroughtDate | PRECTOT | PS | QV2M | T2M | T2MDEW | T2MWET | T2M_MAX |
|---|---|---|---|---|---|---|---|---|---|
| **19300676** | 56043 | 2016-12-28 | 0.02 | 83.33 | 1.41 | -8.71 | -14.10 | -13.84 | -2.49 |
| **19300677** | 56043 | 2016-12-29 | 0.00 | 83.75 | 1.59 | -7.96 | -13.30 | -13.03 | 0.42 |
| **19300678** | 56043 | 2016-12-30 | 1.22 | 82.49 | 2.63 | -2.94 | -7.40 | -7.33 | 3.76 |
| **19300679** | 56043 | 2016-12-31 | 0.44 | 82.19 | 1.75 | -7.56 | -11.98 | -11.82 | -0.95 |

4 rows × 24 columns

```
In [28]:  1  df1 = df.drop(['T2MDEW', 'T2M_MAX', 'T2MWET', 'T2M_MIN', 'WS10M_MAX','WS
          2                 'WS50M', 'WS50M_RANGE', 'QV2M', 'T2M'], axis = 1)
```

In [29]: ►

```
1  plt.figure(figsize=(20,20))
2  corr = df1.corr()
3  kot = corr[corr>=.75]
4  sns.heatmap(kot, cmap="Reds")
```

Out[29]: <AxesSubplot:>



In [12]: ►

```
1  soil_df.head(2)
```

Out[12]:

|   | fips | lat | lon | elevation | slope1 | slope2 | slope3 | slope4 | slope5 | slope6 | ... |
|---|------|-----|-----|-----------|--------|--------|--------|--------|--------|--------|-----|
| 0 | 1001 | 32.536382 | -86.644490 | 63 | 0.0419 | 0.2788 | 0.2984 | 0.2497 | 0.1142 | 0.0170 | ... |
| 1 | 1005 | 31.870670 | -85.405456 | 146 | 0.0158 | 0.1868 | 0.5441 | 0.2424 | 0.0106 | 0.0003 | ... |

2 rows × 32 columns

In [15]: ▶|
```python
1  plt.figure(figsize=(20,20))
2  corr = soil_df.corr()
3  kot = corr[corr>=.8]
4  sns.heatmap(kot, cmap="Reds")
5
```

Out[15]: <AxesSubplot:>



In [20]: ▶|
```python
1  df2 = soil_df.drop(['aspectUnknown', 'CULT_LAND', 'SQ1', 'SQ5', 'SQ6', '
```

In [21]:  ▶|
```python
1  plt.figure(figsize=(20,20))
2  corr = df2.corr()
3  kot = corr[corr>=.75]
4  sns.heatmap(kot, cmap="Reds")
```

Out[21]:  <AxesSubplot:>



In [30]:  ▶|
```python
1  drought_soil = pd.merge(df1, df2, how='left', on=['fips'])
```

In [31]:  ▶|
```python
1  drought_soil = drought_soil[(drought_soil['lat']<= 42) & (drought_soil['
2  drought_soil = drought_soil[(drought_soil['lon']<= -114) & (drought_soil
```

In [32]:  ▶|
```python
1  drought_soil = drought_soil[(drought_soil['DroughtYear'] >=2011) & drough
```

In [33]:  ▶|  `1  drought_soil['DroughtYear'].describe()`

Out[33]:  count    277628.000000
          mean       2015.500411
          std           2.872668
          min        2011.000000
          25%        2013.000000
          50%        2016.000000
          75%        2018.000000
          max        2020.000000
          Name: DroughtYear, dtype: float64

In [34]:  ▶|  `1  drought_soil.to_csv('Data/Drought_Soil.csv')`

In [ ]:  ▶|  `1`

# Appendix A3: Soil and Metereological Data EDA

In [1]:

```python
1   import datetime as dt
2   from pathlib import Path
3   import math
4   import os
5   import sqlite3
6   import json
7   import geopandas as gpd
8   import pygeos
9   import pyproj
10  import shapely
11  import shapely.ops as ops
12  from shapely.geometry import Point, Polygon
13  from shapely.geometry.polygon import Polygon
14  from functools import partial
15
16  import pandas as pd
17  import numpy as np
18  import seaborn as sns
19  import matplotlib.pyplot as plt
20  %matplotlib inline
21
22  from sklearn.model_selection import train_test_split
23
24
25  from sklearn import svm
26  from sklearn.svm import SVC
27  from sklearn.ensemble import RandomForestClassifier
28  from sklearn.naive_bayes import GaussianNB
29  from sklearn.metrics import accuracy_score, classification_report, confus
30
31  from sklearn.feature_selection import SelectKBest
32  from sklearn.feature_selection import chi2, f_classif, mutual_info_class:
33  from functools import partial
34
35
36  from sklearn.preprocessing import StandardScaler
37
38  import warnings
39  warnings.filterwarnings('ignore')
```

## Data Collection: Soil and Metereological

**USA Shape File https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html (https://www.census.gov/geographies/mapping-files/time-series/geo/carto-boundary-file.html)**

In [2]: ▶|
```
1  USA = gpd.read_file("Data/County/cb_2018_us_county_500k.shp")
2  USA.head()
```

Out[2]:

| | STATEFP | COUNTYFP | COUNTYNS | AFFGEOID | GEOID | NAME | LSAD | ALAND |
|---|---|---|---|---|---|---|---|---|
| **0** | 21 | 007 | 00516850 | 0500000US21007 | 21007 | Ballard | 06 | 639387454 |
| **1** | 21 | 017 | 00516855 | 0500000US21017 | 21017 | Bourbon | 06 | 750439351 |
| **2** | 21 | 031 | 00516862 | 0500000US21031 | 21031 | Butler | 06 | 1103571974 |
| **3** | 21 | 065 | 00516879 | 0500000US21065 | 21065 | Estill | 06 | 655509930 |
| **4** | 21 | 069 | 00516881 | 0500000US21069 | 21069 | Fleming | 06 | 902727151 |

◀ ▶

### 3a. Load the dataset 3: US Drought Data - https://power.larc.nasa.gov/ (https://power.larc.nasa.gov/)

In [5]: ▶|
```
1  drought_soil = pd.read_csv('Data/Drought_Soil.csv')
```

### 3b. Daily weather summary data https://www.noaa.gov/ (https://www.noaa.gov/)

In [3]: ▶|

```python
1  # import necessary libraries
2  import glob
3
4  # use glob to get all the csv files
5  # in the folder
6  current_dir = Path(os.getcwd()).absolute()
7  data_dir = current_dir.joinpath('Data')
8  weather_dir = data_dir.joinpath('weather')
9  csv_files = glob.glob(os.path.join(weather_dir, "*.csv"))
10
11
12 # loop over the list of csv files
13 df_list = []
14 for f in csv_files:
15
16     # read the csv file
17     df = pd.read_csv(f)
18     df_list.append(df)
```

In [4]: ▶|

```python
1  ca_daily_df = []
2  for i in df_list:
3      df = pd.DataFrame(i)
4      df = df.dropna()
5      ca_daily_df.append(df)
6  ca_daily_df = pd.concat(ca_daily_df)
```

## Data Preliminary Analysis

In [6]: ▶|

```python
1  # check for missing value
2  def percentMissing(df):
3
4      df_numeric = df.select_dtypes(include=[np.number])
5      numeric_cols = df_numeric.columns.values
6
7      # % of missing data
8      for col in df.columns:
9          # create missing indicator for features with missing data
10         missing = df[col].isnull()
11         pct_missing = np.mean(missing)*100
12         #if pct_missing >60:
13         print('{} - {}%'.format(col, round(pct_missing)))
14         num_missing = np.sum(missing)
```

```
In [7]:     1  # Checking data type
            2  def Datatype(df):
            3      # shape and data types of the data
            4      print("There are {} rows and {} columns".format(df.shape[0], df.shape
            5      print(df.dtypes)
            6
            7      # select numeric columns
            8      df_numeric = df.select_dtypes(include=[np.number])
            9      numeric_cols = df_numeric.columns.values
           10      print(numeric_cols)
           11
           12      # select non numeric columns
           13      df_non_numeric = df.select_dtypes(exclude=[np.number])
           14      non_numeric_cols = df_non_numeric.columns.values
           15      print(non_numeric_cols)
```

## Data Exploration

```
In [8]:     1  drought_soil.head(4)
            2  drought_soil = drought_soil.drop(['Unnamed: 0', 'score'], axis = 1)
```

```
In [9]:     1  drought_soil['DroughtYear'].describe()
```

```
Out[9]:  count    277628.000000
         mean       2015.500411
         std           2.872668
         min        2011.000000
         25%        2013.000000
         50%        2016.000000
         75%        2018.000000
         max        2020.000000
         Name: DroughtYear, dtype: float64
```

```
In [14]:    1  y = drought_soil[drought_soil['lon'] <=-119.5]
            2  z = drought_soil[(drought_soil['lat'] <= 39.5) & (drought_soil['lat'] >=
            3  w = z[(z['lon'] >=-119.5) & (z['lon'] <=-116)]
            4  v = drought_soil[drought_soil['lat'] <= 36.5]
            5  crs = {'init': "EPSG:4326"}
            6  plot_df0 = pd.concat([y,z,w,v])
            7  plot_df0 = plot_df0[~((plot_df0['lat']>39) & (plot_df0['lon']>-120))]
            8  plot_df0 = plot_df0[~((plot_df0['lat']>37.6) & (plot_df0['lon']>-118.5))
            9
           10  geometry = [Point(xy) for xy in zip(plot_df0['lon'], plot_df0['lat'])]
           11  geometry[:3]
           12  geo_soil_df = gpd.GeoDataFrame(plot_df0, crs=crs, geometry=geometry)
```

In [15]: ▶|
```
1  soil_df = geo_soil_df[['fips','lat','lon','elevation', 'slope1','slope2'
2                         'slope4','slope6','slope8','aspectN','aspectE', '
3                         'aspectW','WAT_LAND','NVG_LAND','URB_LAND','GRS_L
4                         'CULTRF_LAND','CULTIR_LAND','SQ2','SQ3', 'SQ4', '
5  soil_df.shape
```

Out[15]:  (361647, 25)

In [16]: ▶|
```
1  soil_df = soil_df[~soil_df.duplicated(['fips'], keep='first')]
2  soil_df.shape
```

Out[16]:  (61, 25)

In [17]: ▶|
```
1  soil_df[['lat','lon','elevation', 'slope1','slope2','slope3','slope4','s
```

Out[17]:

|        | lat        | lon          | elevation    | slope1     | slope2     | slope3     | slope4     | slo      |
|--------|------------|--------------|--------------|------------|------------|------------|------------|----------|
| count  | 61.000000  | 61.000000    | 61.000000    | 61.000000  | 61.000000  | 61.000000  | 61.000000  | 61.00(   |
| mean   | 37.761497  | -120.526784  | 643.426230   | 0.041205   | 0.187852   | 0.117167   | 0.118982   | 0.243    |
| std    | 2.183253   | 2.200278     | 678.960148   | 0.076777   | 0.267800   | 0.112690   | 0.102348   | 0.183    |
| min    | 33.023604  | -123.980998  | 0.000000     | 0.000000   | 0.000000   | 0.000000   | 0.000000   | 0.00(    |
| 25%    | 36.561977  | -122.007205  | 106.000000   | 0.000400   | 0.005300   | 0.029800   | 0.044400   | 0.028    |
| 50%    | 38.021451  | -120.773446  | 444.000000   | 0.001600   | 0.023200   | 0.082700   | 0.096400   | 0.275    |
| 75%    | 39.177739  | -119.749852  | 825.000000   | 0.044500   | 0.322700   | 0.175900   | 0.184900   | 0.402    |
| max    | 41.749903  | -114.038793  | 2630.000000  | 0.343100   | 0.745400   | 0.561200   | 0.483300   | 0.550    |

In [18]: ▶|
```
1  soil_df[['aspectN','aspectE', 'aspectS', 'aspectW']].describe()
```

Out[18]:

|        | aspectN   | aspectE   | aspectS   | aspectW   |
|--------|-----------|-----------|-----------|-----------|
| count  | 61.000000 | 61.000000 | 61.000000 | 61.000000 |
| mean   | 0.151067  | 0.165772  | 0.192116  | 0.230715  |
| std    | 0.093424  | 0.096897  | 0.117666  | 0.137050  |
| min    | 0.000000  | 0.000000  | 0.000000  | 0.000000  |
| 25%    | 0.064200  | 0.089500  | 0.070300  | 0.128400  |
| 50%    | 0.163600  | 0.173000  | 0.201800  | 0.238000  |
| 75%    | 0.240400  | 0.226800  | 0.289000  | 0.317900  |
| max    | 0.348700  | 0.397800  | 0.473800  | 0.560200  |

In [19]: ▶| `1` `soil_df[['WAT_LAND','NVG_LAND','URB_LAND','GRS_LAND', 'FOR_LAND', 'CULTRI`

Out[19]:

| | WAT_LAND | NVG_LAND | URB_LAND | GRS_LAND | FOR_LAND | CULTRF_LAND | CULTIR_ |
|---|---|---|---|---|---|---|---|
| **count** | 61.000000 | 61.000000 | 61.000000 | 61.000000 | 61.000000 | 61.000000 | 61.0 |
| **mean** | 1.033192 | 7.304989 | 2.745081 | 18.985231 | 47.190747 | 1.722027 | 16.1 |
| **std** | 5.806408 | 19.514630 | 13.599561 | 14.832669 | 31.989077 | 2.404207 | 28.7 |
| **min** | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 |
| **25%** | 0.000000 | 0.000000 | 0.000000 | 8.869439 | 18.155117 | 0.000000 | 0.0 |
| **50%** | 0.000000 | 0.000000 | 0.055821 | 15.587294 | 49.677391 | 0.272293 | 0.4 |
| **75%** | 0.000000 | 2.083883 | 0.429140 | 29.065670 | 77.980881 | 3.188647 | 15.7 |
| **max** | 44.035000 | 78.871132 | 99.955193 | 58.796833 | 90.971321 | 9.187908 | 99.9 |

In [20]: ▶| `1` `soil_df[['SQ2','SQ3', 'SQ4']].describe()` *#### these are discrete variable*

Out[20]:

| | SQ2 | SQ3 | SQ4 |
|---|---|---|---|
| **count** | 61.000000 | 61.000000 | 61.000000 |
| **mean** | 1.360656 | 1.573770 | 1.163934 |
| **std** | 1.155174 | 1.257763 | 1.113258 |
| **min** | 0.000000 | 0.000000 | 0.000000 |
| **25%** | 1.000000 | 1.000000 | 1.000000 |
| **50%** | 1.000000 | 1.000000 | 1.000000 |
| **75%** | 1.000000 | 2.000000 | 1.000000 |
| **max** | 7.000000 | 7.000000 | 7.000000 |

In [21]: ▶|
```python
1  fig, ax = plt.subplots(figsize = (10, 10))
2  fig.suptitle('Geospatial Plot of Soil Data (2011-2020)', fontsize=20)
3
4  USA[USA['STATEFP'] == '06'].plot(ax=ax, edgecolor="grey", facecolor="whit
5  soil_df.plot(ax=ax, color="tan", label="Temporal Soil Data")
6  plt.legend(prop={'size':15})
```

Out[21]:  <matplotlib.legend.Legend at 0x1f95eeace20>



Geospatial Plot of Soil Data (2011-2020)

```python
In [22]:
1   # Histograms
2   def histogram(xaxes, df, x, y, nrows, color):
3       plt.rcParams['figure.figsize'] = (x, y)
4
5       fig, axes = plt.subplots(nrows = nrows, ncols = 2)
6       fig.suptitle('Distribution of Meteorological Indicators in West Coast
7
8       # draw histograms in for loop
9       axes = axes.ravel()
10      for idx, ax in enumerate(axes):
11          # drops NaN values
12          ax.hist(df[num_features[idx]].dropna(), bins=40, color= color)
13          ax.set_xlabel(xaxes[idx], fontsize=15)
14          ax.set_ylabel('Counts', fontsize=15)
15          ax.tick_params(axis='both', labelsize=15)
16          right_side = ax.spines["right"]
17          right_side.set_visible(False)
18          top = ax.spines["top"]
19          top.set_visible(False)
20
21      plt.show()
```
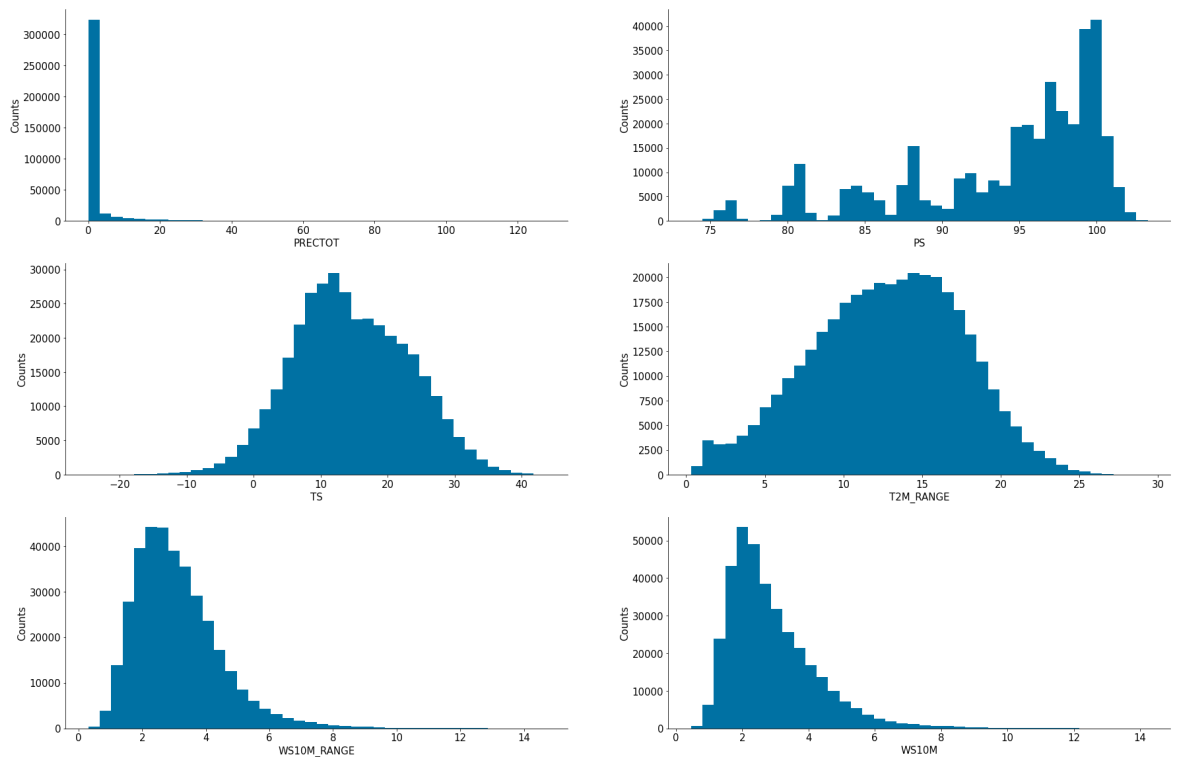
```python
In [23]:
1   # Specify the features of interest
2   num_features = ['PRECTOT', 'PS','TS', 'T2M_RANGE', 'WS10M_RANGE', 'WS10M
3   xaxes = num_features
4   histogram(xaxes, geo_soil_df, 30,20, 3, "#0071A3")
```

Distribution of Meteorological Indicators in West Coast Region (2011-2020)

In [24]:
```python
grouped = geo_soil_df.groupby(by=['DroughtMonth'] ,as_index=False).agg({



df = grouped[["PRECTOT",'PS', 'TS', 'T2M_RANGE','WS10M_RANGE', 'WS10M']]
Month_column = grouped[['DroughtMonth']]
```

In [25]:
```python
sc = StandardScaler()
df_std = sc.fit_transform(df)
df_std = pd.DataFrame(df_std, columns=["PRECTOT",'PS', 'TS', 'T2M_RANGE'
df_std['DroughtMonth'] = Month_column['DroughtMonth']
```

In [26]:
```python
df_std = df_std.replace({'DroughtMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar
                                           5: 'May', 6: 'Jun', 7: 'Jul',
                                           9: 'Sep', 10: 'Oct', 11: 'Nov
fig, ax = plt.subplots(figsize = (10,6))
fig.suptitle('Metereological Indicators: Precipitation vs. TemperatureRan
              x=0.12, y=.95, horizontalalignment='left', verticalalignmen
plot3 = sns.lineplot(data=df_std, x="DroughtMonth", y='T2M_RANGE', color=
                     ax=ax, label="Temp Range at 2 m")
plot4 = sns.lineplot(data=df_std, x="DroughtMonth", y='WS10M', color="#0(
                     ax=ax, label="WindSpeed at 10 m")
plot6 = sns.lineplot(data=df_std, x="DroughtMonth", y="PRECTOT", color="#
                     ax=ax, label="Precipitation")

ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
ax.set_xlabel("Month")
ax.set_ylabel("Scaled Average")

plt.legend(loc="right", bbox_to_anchor=(1.35, 0.5))
sns.despine()
plt.show()
```



Metereological Indicators: Precipitation vs. TemperatureRange vs. WindSpeed

```
In [27]:  ▶| 1  df_std = df_std.replace({'DroughtMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar
             2                                                5: 'May', 6: 'Jun', 7: 'Jul',
             3                                                9: 'Sep', 10: 'Oct', 11: 'Nov
             4  fig, ax = plt.subplots(figsize = (10,6))
             5  fig.suptitle('Metereological Indicators: SP vs. TS vs. WS Range',
             6                x=0.12, y=.95, horizontalalignment='left', verticalalignment
             7  plot1 = sns.lineplot(data=df_std, x="DroughtMonth", y='PS', color="#FFDF0
             8                ax=ax, label="Surface Pressure")
             9  plot2 = sns.lineplot(data=df_std, x="DroughtMonth", y='TS', color="#FFA80
            10                ax=ax, label="Earth Skin Temp")
            11  plot5 = sns.lineplot(data=df_std, x="DroughtMonth", y='WS10M_RANGE', colo
            12                ax=ax, label="WindSpeed Range at 10 m")
            13
            14  ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
            15  ax.set_xlabel("Month")
            16  ax.set_ylabel("Scaled Average")
            17  plt.legend(loc="right", bbox_to_anchor=(1.43, 0.5))
            18  sns.despine()
            19  plt.show()
```



Metereological Indicators: SP vs. TS vs. WS Range

## 1d. Data Preliminary Analysis: Daily Summaries (2010 -2020)

```
In [28]:  ▶| 1  ca_daily_df['DATE'] = pd.to_datetime(ca_daily_df['DATE'])
             2  ca_daily_df['TempYear'] = ca_daily_df['DATE'].dt.year
             3  ca_daily_df['TempMonth'] = ca_daily_df['DATE'].dt.month
             4  ca_daily_df['TempDay'] = ca_daily_df['DATE'].dt.day
```

```
In [29]:  ▶| 1  ca_daily_df = ca_daily_df.rename(columns={'STATION': 'StationCode', 'NAM
             2                                             'LATITUDE': 'StationLatitude',
             3                                             'ELEVATION': 'Elevation', 'DAT
             4                                             'TMAX': 'Max_Temp', 'TMIN': 'M
```

```
In [30]:  ▶| 1  ca_daily_df['Avg_Temp'] = (ca_daily_df['Min_Temp'] + ca_daily_df['Max_Te
             2  ca_daily_df.shape
```

```
Out[30]:  (1314180, 13)
```

```
In [31]:    1  percentMissing(ca_daily_df)
```

```
StationCode - 0%
StationName - 0%
StationLatitude - 3%
StationLongitude - 3%
Elevation - 3%
TempDate - 0%
Precip - 0%
Max_Temp - 0%
Min_Temp - 0%
TempYear - 0%
TempMonth - 0%
TempDay - 0%
Avg_Temp - 0%
```

**Deleting duplicates**

```
In [32]:    1  ca_daily_df.shape
```

Out[32]:  (1314180, 13)

```
In [33]:    1  ca_daily_df = ca_daily_df.sort_values(["Max_Temp", "Min_Temp"], ascending
```

```
In [34]:    1  # check for duplicates in coordinates
            2  ca_daily_df = ca_daily_df[~ca_daily_df.duplicated(['StationLatitude', 'S
            3  ca_daily_df.shape
```

Out[34]:  (1269731, 13)

```
In [35]:    1  grouped1 = ca_daily_df.groupby(by=['TempMonth'] ,as_index=False).agg({'P
            2  df1 = grouped1[["Precip", "Avg_Temp"]]
            3  Month_column1 = grouped1[['TempMonth']]
```

```
In [36]:    1  sc = StandardScaler()
            2  df_std1 = sc.fit_transform(df1)
            3  df_std1 = pd.DataFrame(df_std1, columns=["Precip", "Avg_Temp"])
            4  df_std1['TempMonth'] = Month_column1['TempMonth']
```

```
In [37]: ▶|    1  df_std1 = df_std1.replace({'TempMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar',
               2                                             5: 'May', 6: 'Jun', 7: 'Jul', 8:
               3                                             9: 'Sep', 10: 'Oct', 11: 'Nov',
               4  fig, ax = plt.subplots(figsize = (8,5))
               5  fig.suptitle('Average Monthly: Precipitation vs. Temperature',
               6                x=0.12, y=.95, horizontalalignment='left', verticalalignment
               7  plot1 = sns.lineplot(data=df_std1, x="TempMonth", y='Precip', color="#00
               8                       ax=ax, label="Precipitation")
               9  plot2 = sns.lineplot(data=df_std1, x="TempMonth", y='Avg_Temp', color="#
              10                       ax=ax, label="Temperature")
              11
              12  ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
              13  ax.set_xlabel("Month")
              14  ax.set_ylabel("Scaled Average")
              15  plt.legend(loc="right", bbox_to_anchor=(1.3, 0.5))
              16  sns.despine()
              17  plt.show()
```



```
In [38]: ▶|    1  geometry = [Point(xy) for xy in zip(ca_daily_df['StationLongitude'], ca_
               2  geometry[:3]
               3  geo_daily_df = gpd.GeoDataFrame(ca_daily_df, crs=crs, geometry=geometry)
               4  geo_daily_df.head(2)
```

Out[38]:

| | StationCode | StationName | StationLatitude | StationLongitude | Elevation | TempDate | Pr |
|---|---|---|---|---|---|---|---|
| **99763** | USC00042319 | DEATH VALLEY NATIONAL PARK, CA US | 36.46263 | -116.86720 | -59.1 | 2020-08-16 | |
| **96964** | USW00093115 | IMPERIAL BEACH REAM FIELD NAS, CA US | 32.56797 | -117.11715 | 7.2 | 2010-12-08 | |

```
In [39]:    1  geo_daily_df = geo_daily_df.to_crs({'init': "EPSG:3310"})
```

```
In [40]:    1  geo_daily_df['Precip'].describe()
```

```
Out[40]:  count    1.269731e+06
          mean     6.134752e-02
          std      2.677502e-01
          min      0.000000e+00
          25%      0.000000e+00
          50%      0.000000e+00
          75%      0.000000e+00
          max      1.218000e+01
          Name: Precip, dtype: float64
```

**Plotting Average Monthly Fires with Metereological Indicators**

```
In [41]:    1  grouped2 = geo_fires_df['FireMonth'].value_counts().to_frame()
            2  grouped2 = grouped2.sort_index()
            3  grouped2.reset_index(level=0, inplace=True)
            4  grouped2['AvgFires'] = grouped2['FireMonth'].apply(lambda x: x/10)
            5  grouped2['FireMonth'] = grouped2['index']
            6  grouped2 = grouped2.iloc[: , 1:]
```

```
          ---------------------------------------------------------------------------
          NameError                                 Traceback (most recent call last)
          Input In [41], in <module>
          ----> 1 grouped2 = geo_fires_df['FireMonth'].value_counts().to_frame()
                2 grouped2 = grouped2.sort_index()
                3 grouped2.reset_index(level=0, inplace=True)

          NameError: name 'geo_fires_df' is not defined
```

```
In [ ]:    1  df2 = grouped2[["AvgFires"]]
           2  Month_column2 = grouped2[['FireMonth']]
```

```
In [ ]:    1  sc = StandardScaler()
           2  df_std2 = sc.fit_transform(df2)
           3  df_std2 = pd.DataFrame(df_std2, columns=["AvgFires"])
           4  df_std2['FireMonth'] = Month_column2['FireMonth']
```

```
In [ ]:    1  df_std1 = df_std1.replace({'TempMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar'
           2                                            5: 'May', 6: 'Jun', 7: 'Jul', 8:
           3                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
           4  df_std2 = df_std2.replace({'FireMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar'
           5                                            5: 'May', 6: 'Jun', 7: 'Jul', 8:
           6                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
           7  fig, ax = plt.subplots(figsize = (8, 5))
           8  fig.suptitle('Average Monthly: Precipitation vs. Fires', fontsize=16)
           9  plot1 = sns.lineplot(data=df_std1, x="TempMonth", y='Precip', color="#00
          10                       ax=ax, label="Precipitation")
          11  plot4 = sns.lineplot(data=df_std2, x="FireMonth", y='AvgFires', color="#
          12                       ax=ax, label="Fires")
          13  ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
          14  ax.set_xlabel("Month")
          15  ax.set_ylabel("Scaled Average")
          16  plt.legend(loc="upper right")
          17  sns.despine()
          18  plt.show()
```

```
In [ ]:    1  df_std1 = df_std1.replace({'TempMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar'
           2                                            5: 'May', 6: 'Jun', 7: 'Jul', 8:
           3                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
           4  df_std2 = df_std2.replace({'FireMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar'
           5                                            5: 'May', 6: 'Jun', 7: 'Jul', 8:
           6                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
           7  fig, ax = plt.subplots(figsize = (8,5))
           8  fig.suptitle('Average Monthly: Temperatures vs. Fires', fontsize=16)
           9  plot2 = sns.lineplot(data=df_std1, x="TempMonth", y='Avg_Temp', color="#
          10                       ax=ax, label="Temperature")
          11  plot4 = sns.lineplot(data=df_std2, x="FireMonth", y='AvgFires', color="#
          12                       ax=ax, label="Fires")
          13  ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
          14  ax.set_xlabel("Month")
          15  ax.set_ylabel("Scaled Average")
          16  plt.legend(loc="upper right")
          17  sns.despine()
          18  plt.show()
```

```
In [ ]:    1  df_std2 = df_std2.replace({'FireMonth' : {1: 'Jan', 2 : 'Feb', 3 : 'Mar'
           2                                            5: 'May', 6: 'Jun', 7: 'Jul', 8:
           3                                            9: 'Sep', 10: 'Oct', 11: 'Nov',
           4  fig, ax = plt.subplots(figsize = (8, 5))
           5  fig.suptitle('Average Monthly: WindSpeed vs. Fires', fontsize=16)
           6  plot3 = sns.lineplot(data=df_std, x="DroughtMonth", y='WS10M', color="#3
           7                       ax=ax, label="WindSpeed")
           8  plot4 = sns.lineplot(data=df_std2, x="FireMonth", y='AvgFires', color="#
           9                       ax=ax, label="Fires")
          10  ax.set_xticks([0,1,2,3,4,5,6,7, 8,9, 10,11])
          11  ax.set_xlabel("Month")
          12  ax.set_ylabel("Scaled Average")
          13  plt.legend(loc="upper right")
          14  sns.despine()
          15  plt.show()
```

In [ ]:

```python
r_df = df_std1[['Precip', 'Avg_Temp']]
r_df['WS'] = df_std['WS10M']
r_df['Avg_Fires'] = df_std2["AvgFires"]
```

In [ ]:

```python
corrmat = r_df.corr()
top = corrmat.index
plt.figure(figsize=(8,8))

            #plot heat map
g=sns.heatmap(r_df[top].corr(),annot=True,cmap="Reds")
```

In [ ]:

```python

```

# Appendix A4 - Data Merging Code

```
In [1]:    1  import datetime as dt
           2  from pathlib import Path
           3  import math
           4  import os
           5  import sqlite3
           6  import json
           7  import geopandas as gpd
           8  import pygeos
           9  import pyproj
          10  import shapely
          11  import shapely.ops as ops
          12  from shapely.geometry import Point, Polygon
          13  from shapely.geometry.polygon import Polygon
          14  from functools import partial
          15
          16  import pandas as pd
          17  import numpy as np
          18  import seaborn as sns
          19  import matplotlib.pyplot as plt
          20  %matplotlib inline
          21
          22  from sklearn.model_selection import train_test_split
          23
          24
          25  from sklearn import svm
          26  from sklearn.svm import SVC
          27  from sklearn.ensemble import RandomForestClassifier
          28  from sklearn.naive_bayes import GaussianNB
          29  from sklearn.metrics import accuracy_score, classification_report, confu
          30
          31  from sklearn.feature_selection import SelectKBest
          32  from sklearn.feature_selection import chi2, f_classif, mutual_info_class
          33  from functools import partial
          34
          35
          36  from sklearn.preprocessing import StandardScaler
          37
          38  import warnings
          39  warnings.filterwarnings('ignore')
```

## Dataset 1

```
In [250]:  ▶  1  print(geo_daily_df.crs)
              2  print(geo_nasa_df.crs)
              3  print(geo_fires_df.crs)
              4  print(geo_soil_df.crs)
```

```
+init=epsg:3310 +type=crs
+init=epsg:3310 +type=crs
+init=epsg:3310 +type=crs
+init=epsg:4326 +type=crs
```

**First combine all weather and Soil data with nasa dataset**

```
In [251]:  ▶  1  def get_nearestpoint(df1, df1day, df2, df2day, days, dist):
              2      """
              3          This Function merges dataframe for selected day by finding neares
              4          for each day and creates mini dfs for each day of month
              5      """
              6
              7      dfs = []
              8      for day in days:
              9          df = df1[df1[df1day] == day]
             10          df3 = df2[df2[df2day] == day]
             11          m_df = gpd.sjoin_nearest(df, df3, how='left', distance_col=dist)
             12          m_df[dist] = m_df[dist].apply(lambda x: x/1000)
             13          d = pd.DataFrame(m_df)
             14          dfs.append(d)
             15
             16      dfs = pd.concat(dfs)
             17      return dfs
```

```
In [252]:  ▶  1  def merge_data(data1, df1year, df1month, df1day, data2, df2year, df2montl
              2      """
              3          This Function filters dataframe by year and months and calls for
              4          append it and then converts it into pandas df.
              5      """
              6      dfs = []
              7      for month in months:
              8          df1 = data1[(data1[df1year] == year) & (data1[df1month] == month
              9          df2 = data2[(data2[df2year] == year) & (data2[df2month] == month
             10          df = get_nearestpoint(df1, df1day, df2, df2day, days, dist)
             11          dfs.append(df)
             12
             13      dfs = pd.concat(dfs)
             14      return dfs
```

```
In [253]:  ▶|  1  def get_data(df1, df1year, df1month,df1day, df2, df2year, df2month,df2day
              2      """
              3          This calls for all dataframes and combine it and create one datas
              4          nasa data and daily temperatures
              5      """
              6      years = list(range(2010, 2021))
              7      months = list(range(1, 13))
              8      days = list(range(1, 32))
              9
             10      dataframesList = []
             11      for year in years:
             12          data = merge_data(df1, df1year, df1month,df1day, df2, df2year, d
             13
             14          dataframesList.append(data)
             15
             16      df = gpd.GeoDataFrame(pd.concat(dataframesList), crs=crs)
             17      try:
             18          df.drop('index_right', axis=1, inplace=True)
             19      except ValueError:
             20          # ignore if there are no index columns
             21          pass
             22
             23      print(df.shape)
             24
             25      return df
             26
```

**Combine California fire data with nasa DF**

```
In [254]:  ▶|  1  pixel_temp = get_data(geo_nasa_df, 'ActiveYear', 'ActiveMonth', 'ActiveDa
              2                        geo_daily_df, 'TempYear', 'TempMonth', 'TempDay',
```

(114599, 31)

```
In [255]:  ▶|  1  pixel_soil = get_data(pixel_temp, 'ActiveYear', 'ActiveMonth', 'ActiveDay
              2                        geo_soil_df, 'DroughtYear', 'DroughtMonth', 'Drough
```

(119665, 66)

**Dropping Duplicates**

```
In [256]:  ▶|  1  pixel_soil = pixel_soil[~(pixel_soil.index.duplicated(keep='first'))]
```

```
In [257]:  ▶|  1  pixel_fire= get_data(pixel_soil, 'ActiveYear', 'ActiveMonth', 'ActiveDay
              2                        geo_fires_df, 'FireYear','FireMonth', 'FireDay', 'f
```

(114599, 76)

```
In [258]:  ▶|  1  merged_df =pixel_fire
```

```
In [259]:   ▶  1  # create a set for important variables and target variables.
               2  # Delete repetitive dates, Station name, station codes are not needed.
               3  # Name of the fire is also irrelevant and acres are also not
               4  # needed because they are attributes that are logged after the fire event
               5  merged_df = merged_df.drop(['Satellite', 'StationCode', 'StationName', '
               6                              'TempDate', 'TempYear', 'TempMonth', 'TempDay
               7                              'DroughtYear', 'DroughtMonth', 'DroughtDay',
               8                              'FireYear', 'Name', 'UniqueId'], axis = 1)
```

**Perform Spatial Analysis**

```
In [260]:   ▶  1  labeled_data = merged_df[merged_df['fire_dist'].notnull()]
               2  unlabeled_data = merged_df[merged_df['fire_dist'].isnull()] # taking pixe
```

```
In [261]:   ▶  1  print(labeled_data.shape)
               2  print(unlabeled_data.shape)
```

```
(82493, 57)
(32106, 57)
```

```
In [262]:   ▶  1  labeled_data[['fire_dist', 'Station_dist', 'Drought_dist']].describe()
```

Out[262]:

|       | fire_dist | Station_dist | Drought_dist |
|-------|-----------|--------------|--------------|
| count | 82493.000000 | 82493.000000 | 82493.000000 |
| mean | 299.182305 | 50.681572 | 310.481150 |
| std | 230.135528 | 85.893810 | 158.940170 |
| min | 0.000000 | 0.051470 | 0.260367 |
| 25% | 122.804595 | 14.653380 | 202.725861 |
| 50% | 237.723699 | 23.041096 | 299.132089 |
| 75% | 430.907748 | 36.686614 | 442.452498 |
| max | 1259.457214 | 518.527597 | 805.499278 |

The maximum number of TotalAcres burned is 410202 Acres and roughly 1660 km^2, meaning the maximum distance a potential fire goes in any one direction is approximately around 830 km and min 1 km. This could've been a threshold to filter our data out for fire_dist to fire pixel if we were only trying to detect the fire pixel is true fire or not in any given day of the fire (considering fire can go on for weeks). However, the main purpose of this project is to build a model that can be used for early detection of the wildfire, so the hazard can be prevented from spreading. Not all fire pixels are true fire pixels and some times it is a false alarm. So we will filter for pixels that are within than 1 km from the true fire event, 1 km is used as a threshold because MODIS location coordinates are center of 1km fire pixel but not necessarily the actual location of the fire as one or more fires can be detected within the 1km pixel. Randomly sampled data from unlabeled data will be used as false alarms because no fire event was mapped to those dates.

In [263]: ▶| 
```
1 class1 = labeled_data[labeled_data['fire_dist'] <=1]
```

In [264]: ▶| 
```
1 class1[['fire_dist', 'Station_dist']].describe()
```

Out[264]:

|       | fire_dist | Station_dist |
|-------|-----------|--------------|
| count | 2087.000000 | 2087.000000 |
| mean  | 0.218470 | 20.609494 |
| std   | 0.285376 | 12.244417 |
| min   | 0.000000 | 0.177639 |
| 25%   | 0.000000 | 11.653871 |
| 50%   | 0.030367 | 19.253845 |
| 75%   | 0.410110 | 27.747893 |
| max   | 0.999613 | 71.508454 |

Station distance is sort of irrelevant as it is maximum 71 km away from the fire pixel, which is usually within county limit, and drastic weather changes are highly unlikely for such close approximation.

In [265]: ▶| 
```
1 # class2 = unlabeled_data.sample(frac=.08)
2 unlabeled_fires = labeled_data[labeled_data['fire_dist'] >1]
```

In [957]: ▶| 
```
1 #class2 = class1.append(class2)
```

In [266]: ▶| 
```
1 print(class1.shape)
2 print(unlabeled_fires.shape)
3 print(unlabeled_data.shape)
```

```
(2087, 57)
(80406, 57)
(32106, 57)
```

**Check for duplicates after merging**

In [267]: ▶| 
```
1 class1[class1.index.duplicated(keep=False)]
```

Out[267]:

| NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Confidence | BrightT31 | |
|--------------|---------------|------------|------|-------|------------|------------|-----------|---|

0 rows × 57 columns

◄ ▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮▮                                                ▶

In [268]: ▶|  `1  unlabeled_fires[unlabeled_fires.index.duplicated(keep=False)]`

Out[268]:

| NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Confidence | BrightT31 |
|---|---|---|---|---|---|---|---|

0 rows × 57 columns

◀ | ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓ | | ▶

In [269]: ▶|  `1  unlabeled_data[unlabeled_data.index.duplicated(keep=False)]`

Out[269]:

| NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Confidence | BrightT31 |
|---|---|---|---|---|---|---|---|

0 rows × 57 columns

◀ | ▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓▓ | | ▶

In [270]: ▶|
```
1  class1["Target"] = 1
2  unlabeled_fires["Target"] = 0
3  unlabeled_data["Target"] = 0
4  class2 = unlabeled_fires.append(unlabeled_data)
```

In [273]: ▶|  `1  fire_data = class1.append(class2)`

In [274]: ▶|  `1  fire_data.shape`

Out[274]:  (114599, 58)

In [275]: ▶|  `1  fire_data = fire_data.replace({'DayNight': {'D':1, 'N':0}, 'ConfidenceBi`

In [276]: ▶|
```
1  features = fire_data.drop(['ActiveYear', 'ActiveMonth', 'ActiveDay','Sta
2                           'lat','lon','Drought_dist', 'FireDay', 'fire_
3                           'ActiveDate', 'Confidence', 'geometry', 'Fire
4  features.shape
```

Out[276]:  (114599, 43)

In [277]: ▶|  `1  fire_data.to_csv('Data/clean_dataset_preliminary.csv')`

## Final Analysis: Secondary and Optional Approach for Modeling

For this analysis, we will be going back to data preparation part and change the data little bit. For the first approach we took out all the noise data and trained the model with fire pixels that were mapped to the true fire events within the 1 km of the fire event on the day of fire alarm and with the unlabeled fire pixels that were not mapped with the fire events. However, this approach has some biases:

1. We do not know how much of the unlabeled data (fire pixels that was not mapped to the exact date of fire event) are true fire pixels after the start of fire.

2. All the fire pixels that were about more than 1 km away could also be true fire events on a different day for different fire event, because fire can go on for weeks. For example, for class C fire which is less than 100 acres (0.400639 km^2) can go minimum 0.400639 km and maximum 1 km distance for fire, any fire pixels in that range of distance and duration range between the day of actual fire and the fire containment date can be considered true fire event, anything outside of that distance threshold can be a false alarm, however anything outside of duration but within the 1 km distance means it is an ongoing fire. Note: The containment date signifies that control line has been completed around the fire, and any associated spot fires, which can reasonably be expected to stop the fire's spread, but the fire can continue going for months.

There are few assumptions I will be making when creating a transformed dataset only for this project purpose.
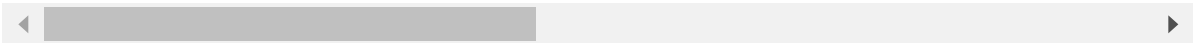
## Dataset 2

In [278]: ▶ | 1  `class1.head() ## True Fire pixels on the day of fire and with distance w`

Out[278]:

|  | NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Confidence | Bright |
|---|---|---|---|---|---|---|---|---|
| **39687** | 36.8878 | -118.2121 | 315.9 | 1.3 | 1.1 | 2011-03-02 | 74 | 28 |
| **39705** | 36.8816 | -118.2201 | 305.3 | 1.5 | 1.2 | 2011-03-02 | 50 | 29 |
| **77188** | 37.3252 | -118.5696 | 330.4 | 3.8 | 1.8 | 2011-05-25 | 0 | 28 |
| **77769** | 35.6501 | -118.3755 | 341.0 | 2.8 | 1.6 | 2011-05-27 | 88 | 29 |
| **77770** | 35.6372 | -118.3795 | 320.3 | 2.8 | 1.6 | 2011-05-27 | 40 | 30 |

5 rows × 58 columns

In [279]: ▶|

```python
def get_firepoints(data1, df1year, df1month, df1day, data2, df2year, df2m
    """
        This Function filters dataframe by year and months and calls for
        append it and then converts it into pandas df.
    """
    months = list(range(1, 13))
    dfs = []
    for month in months:
        df1 = data1[(data1[df1year] == year) & (data1[df1month] == month
        df2 = data2[(data2[df2year] == year) & (data2[df2month] == month

        df = gpd.sjoin_nearest(df1, df2, how='left', distance_col=dist)
        df[dist] = df[dist].apply(lambda x: x/1000)
        d = pd.DataFrame(df)
        dfs.append(d)

    dfs = pd.concat(dfs)
    return dfs
```
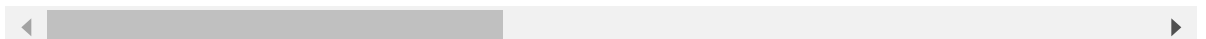
In [280]: ▶|

```python
def get_fire_data(df1, df1year, df1month, df1day, df2, df2year, df2month
    """
        This calls_ for all dataframes and combine it and create one data
        fire data, so we can use the combined information to find the est
    """
    years = list(range(2011, 2021))
    dataframesList = []
    for year in years:
        data = get_firepoints(df1, df1year, df1month, df1day, df2, df2yea
        dataframesList.append(data)

    df = gpd.GeoDataFrame(pd.concat(dataframesList), crs=crs)
    try:
        df.drop('index_right', axis=1, inplace=True)
    except ValueError:
        # ignore if there are no index columns
        pass

    print(df.shape)

    return df
```

In [281]: ▶|

```python
pixel_soil.head(2)
```

Out[281]:

| | NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Satellite | Confidenc |
|---|---|---|---|---|---|---|---|---|
| **26628** | 34.6033 | -118.3848 | 305.8 | 1.2 | 1.1 | 2011-01-06 | Terra | 5 |
| **26638** | 39.8467 | -121.5210 | 308.2 | 1.1 | 1.0 | 2011-01-07 | Terra | 7 |

2 rows × 66 columns

In [282]: ▶|

```python
geo_fires_df = geo_fires_df[geo_fires_df['TotalAcres'].notnull()]
```

In [283]: ▶|
```
1  geo_fires_df['TotalAcres_sq_km'] = geo_fires_df['TotalAcres'].apply(lamb
```

In [284]: ▶|
```
1  print(geo_fires_df.shape)
2  geo_fires_df.head(2)
```

(3795, 11)

Out[284]:

|   | FireCause | TotalAcres | geometry | FireDate | FireYear | FireMonth | FireDay | Name | U |
|---|-----------|------------|----------|----------|----------|-----------|---------|------|---|
| **0** | Powerline | 109.60250 | MULTIPOLYGON (((-116842.172 97942.739, -116837... | 2020-06-18 | 2020.0 | 6.0 | 18.0 | nelson | |
| **1** | Equipment Use | 685.58502 | MULTIPOLYGON (((-117329.343 90212.620, -117322... | 2020-06-01 | 2020.0 | 6.0 | 1.0 | amoruso | |

◀                                         ▶

In [285]: ▶|
```
1  pixel_fire2= get_fire_data(pixel_soil, 'ActiveYear', 'ActiveMonth', 'Act
2                             geo_fires_df, 'FireYear','FireMonth', 'FireDay
```

(114599, 77)

In [286]: ▶|
```
1  merged_df2 = pixel_fire2
```

In [287]: ▶|
```
1  # create a set for important variables and target variables.
2  # Delete repetitive dates, Station name, station codes are not needed.
3  # Name of the fire is also irrelevant and acres are also not
4  # needed because they are attributes that are logged after the fire even
5  merged_df2 = merged_df2.drop(['Satellite', 'StationCode', 'StationName',
6                               'TempDate', 'TempYear', 'TempMonth', 'TempDay
7                               'DroughtYear', 'DroughtMonth', 'DroughtDay',
```

In [288]: ▶|
```
1  mapped_fire = merged_df2[merged_df2['fire_dist'].notnull()]
2  unmapped_fire = merged_df2[merged_df2['fire_dist'].isnull()] # taking pi
```

In [289]: ▶|
```
1  print(mapped_fire.shape)
2  print(unmapped_fire.shape)
```

(113912, 61)
(687, 61)

In [290]: ▶|
```python
1  mapped_fire['fire_dist'].describe()
```

Out[290]:
```
count    113912.000000
mean         79.608633
std         127.814990
min           0.000000
25%           0.000000
50%          29.523428
75%          87.483079
max        1151.519722
Name: fire_dist, dtype: float64
```

In [291]: ▶|
```python
1  def get_duration(df):
2      df['Active_minus_FireDate'] = (df["ActiveDate"] - df["FireDate"]).dt
3      df['Area_diff'] = (df["TotalAcres_sq_km"] - df["fire_dist"])
4      df = df.sort_values('Active_minus_FireDate', ascending=True)
5
6      return df
```

In [292]: ▶|
```python
1  mapped_fire = get_duration(mapped_fire)
```

In [293]: ▶|
```python
1  # fire distance bigger than total Acres means fire is outside of range a
2  # We used 32 km (10 miles) as a threshold for the area. Any fire outside
3  all_false = mapped_fire[(mapped_fire['TotalAcres_sq_km'] <=100) &
4                          (mapped_fire['fire_dist'] >100) & (mapped_fire['
5
```

In [294]: ▶|
```python
1  all_fire = mapped_fire[~((mapped_fire['TotalAcres_sq_km'] <=100) &
2                           (mapped_fire['fire_dist'] >100) &
3                           (mapped_fire['Area_diff'] <0))].sort_values('fi
```

In [295]: ▶|
```python
1  f = all_fire[(all_fire['TotalAcres_sq_km'] <200) &
2               (all_fire['fire_dist'] >200) &
3               (all_fire['Area_diff'] <0)]
4
5  all_false = all_false.append(f)
```

In [296]: ▶|
```python
1  all_fire = all_fire[~((all_fire['TotalAcres_sq_km'] <200) &
2                        (all_fire['fire_dist'] >200) &
3                        (all_fire['Area_diff'] <0))]
```

In [297]: ▶|
```python
1  f = all_fire[(all_fire['TotalAcres_sq_km'] <400) &
2               (all_fire['fire_dist'] >400) &
3               (all_fire['Area_diff'] <0)]
4
5  all_false = all_false.append(f)
```

In [298]: ▶|
```python
1  all_fire = all_fire[~((all_fire['TotalAcres_sq_km'] <400) &
2                        (all_fire['fire_dist'] >400) &
3                        (all_fire['Area_diff'] <0))]
```

In [299]: ▶|
```python
1  f = all_fire[(all_fire['TotalAcres_sq_km'] <5) &
2               (all_fire['fire_dist'] >5.99) &
3               (all_fire['Area_diff'] <0)]
4
5  all_false = all_false.append(f)
```

In [300]: ▶|
```python
1  all_fire = all_fire[~((all_fire['TotalAcres_sq_km'] <5) &
2                        (all_fire['fire_dist'] >5.99) &
3                        (all_fire['Area_diff'] <0))]
```

In [301]: ▶|
```python
1  f = all_fire[(all_fire['Area_diff'] <0)]
2  all_false = all_false.append(f)
3  all_false = all_false.append(unmapped_fire)
```

In [302]: ▶|
```python
1  all_fire = all_fire[~(all_fire['Area_diff'] <0)]
```

### Removing duplicates from all the merging

In [303]: ▶|
```python
1  False_pixels = all_false[~all_false.index.duplicated(keep='first')]
2  False_pixels = get_duration(False_pixels)
3  False_pixels.shape
```

Out[303]: (67863, 63)

In [304]: ▶|
```python
1  True_pixels = all_fire[~all_fire.index.duplicated(keep='first')]
2  True_pixels = get_duration(True_pixels)
3  True_pixels.shape
```

Out[304]: (46736, 63)

### Filtering for California Pixels only

In [305]: ▶|
```python
1  x = False_pixels[~((False_pixels['NasaLongitude'] >-119.8) & (False_pixel
2  y = x[~((x['NasaLongitude'] >-119) & (x['NasaLatitude'] >38))]
3  False_pixels = y[~((y['NasaLongitude'] >-118) & (y['NasaLatitude'] >35.9]
4  False_pixels.shape
```

Out[305]: (57935, 63)

In [306]: ▶|
```python
1  geometry1 = [Point(xy) for xy in zip(False_pixels['NasaLongitude'], False
2  geometry1[:3]
3  plot_df1 = gpd.GeoDataFrame(False_pixels, crs=crs, geometry=geometry1)
```

In [307]: ▶|
```python
1  True_pixels = True_pixels[~((True_pixels['NasaLongitude'] >-119.8) & (Tru
2  True_pixels.shape
```

Out[307]: (45186, 63)

```
In [308]:   ▶  1  geometry2 = [Point(xy) for xy in zip(True_pixels['NasaLongitude'], True_p
               2  geometry2[:3]
               3  plot_df2 = gpd.GeoDataFrame(True_pixels, crs=crs, geometry=geometry2)
```
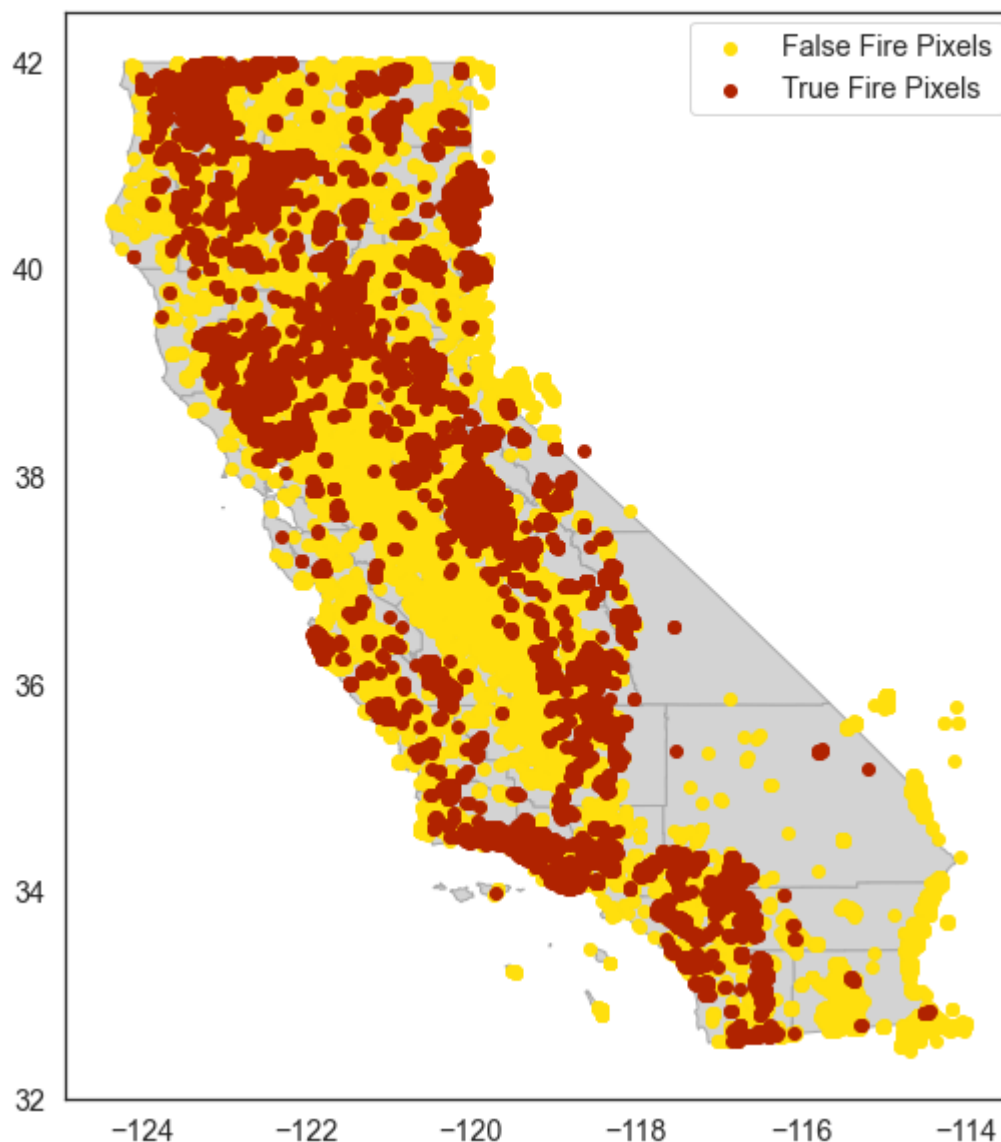
**Plotting all Labeled and Unlabeled Fires**

In [309]: ▶|

```python
1  fig, ax = plt.subplots(figsize = (10,10))
2  fig.suptitle('Mapped All Fire Pixels: labeled and unlabeled (2011-2020)'
3  plt.yticks([32, 34, 36, 38, 40, 42])
4  plt.xticks([-124, -122, -120, -118, -116, -114])
5
6  plt0 = USA[USA.STATEFP == '06'].plot(ax = ax, edgecolor="darkgrey", face
7  plt1 = plot_df1.plot(ax=ax, color="#FFDF0D", label="False Fire Pixels")
8  plt2 = plot_df2.plot(ax=ax, color="#B02300", label="True Fire Pixels")
9
10 handles, labels = ax.get_legend_handles_labels()
11 fig.legend(handles, labels, loc=(0.65,0.8))
12 plt.show()
```

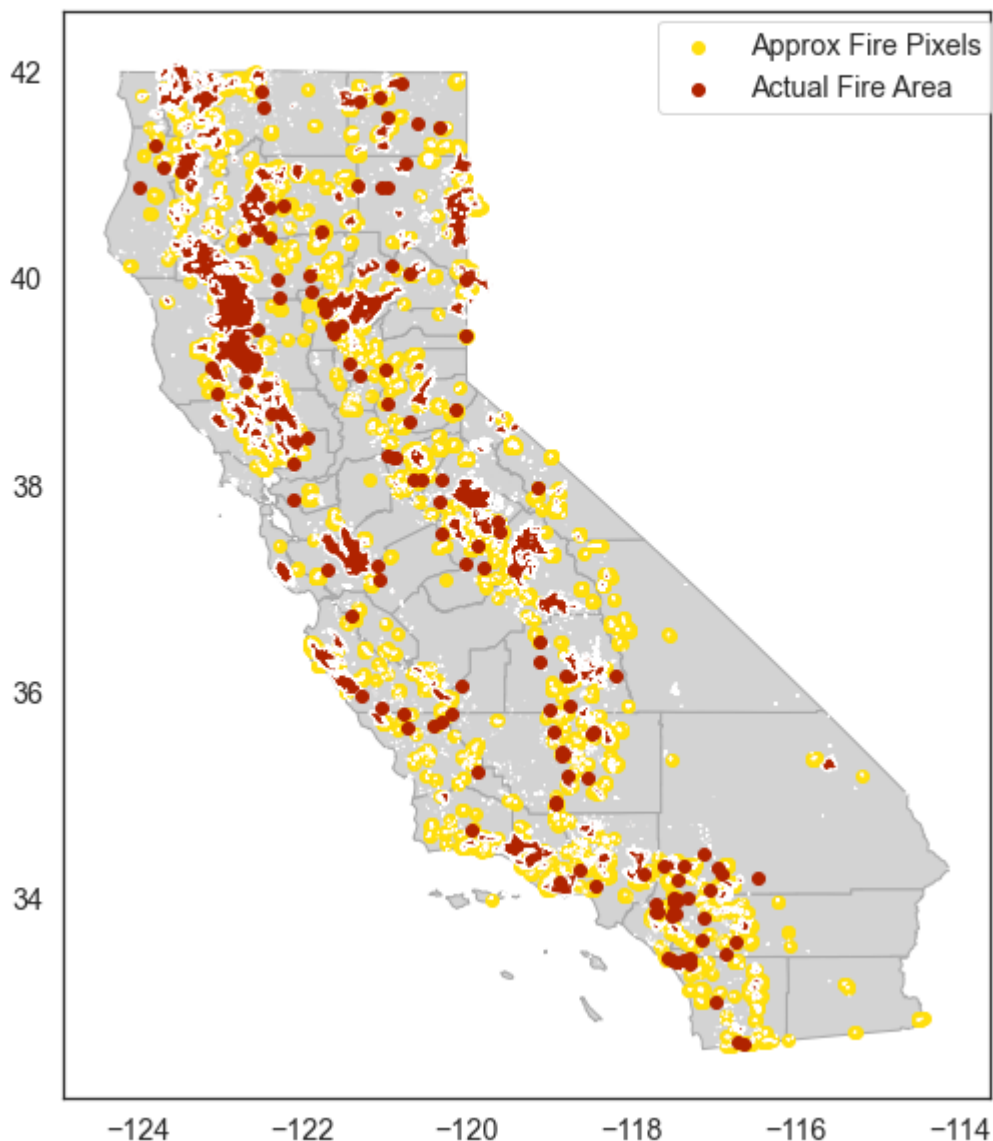Mapped All Fire Pixels: labeled and unlabeled (2011-2020)



In [310]: ▶|

```python
1  plot_gf1 = geo_fires_df.to_crs({'init': "EPSG:4326"})
```

In [311]: ▶|
```
1  plot_df2 = plot_df2[(plot_df2['Active_minus_FireDate'] >=-1) & (plot_df2
```

In [312]: ▶|
```
1  fig, ax = plt.subplots(figsize = (10,10))
2  fig.suptitle('Mapped Labeled Actual Fires with Fire Pixel (2011-2020)',
3  plt.yticks([32, 34, 36, 38, 40, 42])
4  plt.xticks([-124, -122, -120, -118, -116, -114])
5
6  plt0 = USA[USA.STATEFP == '06'].plot(ax = ax, edgecolor="darkgrey", face
7  plt1 = plot_df2.plot(ax=ax, color="#FFDF0D", label="Approx Fire Pixels")
8  plt2 = plot_gf1.plot(ax=ax, color="#B02300", label ="Actual Fire Area")
9  handles, labels = ax.get_legend_handles_labels()
10 fig.legend(handles, labels, loc=(0.63,0.8))
11
12 plt.show()
```



Mapped Labeled Actual Fires with Fire Pixel (2011-2020)

```
In [313]:    1  dta0 = False_pixels.drop(['fire_dist','TotalAcres_sq_km', 'Active_minus_
             2
             3  dta1 = True_pixels.drop(['fire_dist','TotalAcres_sq_km', 'Active_minus_F:
```

```
In [314]:    1  print(dta0.shape)
             2  print(dta1.shape)
```

```
(57935, 59)
(45186, 59)
```

```
In [315]:    1  dta1["Target"] = 1
             2  dta0["Target"] = 0
```

```
In [316]:    1  new_data = dta1.append(dta0)
```

```
In [317]:    1  new_data = new_data.replace({'DayNight': {'D':1, 'N':0}})
             2  new_data = new_data.reset_index(drop=True)
```

```
In [318]:    1  new_data = new_data.sort_values(['ActiveDate'])
             2
             3  new_data['geometry'] = list(zip(new_data['NasaLongitude'], new_data['Nas:
```

```
In [319]:    1  new_df = new_data.set_index(['ActiveDate', 'geometry'])
```

```
In [320]:    1  new_df.to_csv('Data/clean_dataset1.csv')
```

```
In [ ]:      1
```

# Appendix A.5 Data Modeling Preliminary - Random Forest and SVM - Classification Model

In [1]:

```python
import datetime as dt
from pathlib import Path
import math
import os
import json

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.model_selection import train_test_split


from sklearn import svm
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, classification_report, confu

from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2, f_classif, mutual_info_class
from functools import partial

from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedKFold
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_
from scipy import stats
from sklearn.preprocessing import StandardScaler

# Import Keras
from keras.models import Sequential
from keras.layers import Dense, LSTM, BatchNormalization
from keras.optimizers import Adam
from keras.callbacks import ReduceLROnPlateau, EarlyStopping
from keras.regularizers import l2
from time import time
import pickle

import warnings
warnings.filterwarnings('ignore')
```

```
Using TensorFlow backend.
```

In [2]: ▶|
```python
1  fires_data = pd.read_csv("Data/clean_dataset_preliminary.csv")
2  fires_data.drop(columns=fires_data.columns[0], axis=1, inplace=True)
3  fires_data.head(2)
```

Out[2]:

| | NasaLatitude | NasaLongitude | Brightness | Scan | Track | ActiveDate | Confidence | BrightT31 |
|---|---|---|---|---|---|---|---|---|
| **0** | 36.8878 | -118.2121 | 315.9 | 1.3 | 1.1 | 2011-03-02 | 74 | 280.6 |
| **1** | 36.8816 | -118.2201 | 305.3 | 1.5 | 1.2 | 2011-03-02 | 50 | 291.6 |

2 rows × 58 columns

◀ |▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬| ▶

In [4]: ▶|
```python
1  features = fires_data.drop(['ActiveYear', 'ActiveMonth', 'ActiveDay','St:
2                             'lat','lon','Drought_dist', 'FireDay', 'fire_(
3                             'ActiveDate', 'ConfidenceBinned', 'geometry',
4  features.shape
```

Out[4]: (114599, 42)

In [5]: ▶|
```python
1  target = fires_data['Target']
2  target.shape
```

Out[5]: (114599,)

In [6]: ▶|
```python
1  # Checking data type
2  def Datatype(df):
3      # shape and data types of the data
4      print("There are {} rows and {} columns".format(df.shape[0], df.shape
5      print(df.dtypes)
6
7      # select numeric columns
8      df_numeric = df.select_dtypes(include=[np.number])
9      numeric_cols = df_numeric.columns.values
10     print(numeric_cols)
11
12     # select non numeric columns
13     df_non_numeric = df.select_dtypes(exclude=[np.number])
14     non_numeric_cols = df_non_numeric.columns.values
15     print(non_numeric_cols)
```

In [7]: ▶|
```python
1  cor_matrix = features.corr().abs()
2  upper_tri = cor_matrix.where(np.triu(np.ones(cor_matrix.shape),k=1).asty|
3
4  to_drop = [column for column in upper_tri.columns if any(upper_tri[colum|
5  print(to_drop)
```

```
['Track', 'Avg_Temp', 'elevation', 'slope2', 'slope3', 'slope6', 'slope8',
'aspectN', 'aspectE', 'aspectS', 'aspectW', 'NVG_LAND', 'FOR_LAND', 'CULTRF
_LAND', 'CULTIR_LAND', 'SQ2', 'SQ3']
```

In [8]: ▶|   1   `features = features.drop(to_drop, axis = 1)`

In [10]: ▶|   1   `Datatype(features)`

```
There are 114599 rows and 25 columns
NasaLatitude     float64
NasaLongitude    float64
Brightness       float64
Scan             float64
Confidence         int64
BrightT31        float64
Frp              float64
DayNight           int64
HotSpotType      float64
Elevation        float64
Precip           float64
Max_Temp         float64
Min_Temp         float64
PRECTOT          float64
PS               float64
T2M_RANGE        float64
TS               float64
WS10M            float64
WS10M_RANGE      float64
slope1           float64
slope4           float64
WAT_LAND         float64
URB_LAND         float64
GRS_LAND         float64
SQ4                int64
dtype: object
['NasaLatitude' 'NasaLongitude' 'Brightness' 'Scan' 'Confidence'
 'BrightT31' 'Frp' 'DayNight' 'HotSpotType' 'Elevation' 'Precip'
 'Max_Temp' 'Min_Temp' 'PRECTOT' 'PS' 'T2M_RANGE' 'TS' 'WS10M'
 'WS10M_RANGE' 'slope1' 'slope4' 'WAT_LAND' 'URB_LAND' 'GRS_LAND' 'SQ4']
[]
```

## Preliminary Modeling

In [11]:

```python
def results(classifier, x_train, y_train, x_test, y_test):

    model = classifier.fit(x_train, y_train)

    y_pred = classifier.predict(x_test)

    #Checking the accuracy
    accuracy = round(accuracy_score(y_test, y_pred)*100,2)
    print("Accuracy score: {}{}".format(round(accuracy, 2), '%'))

    cv_accuracy_score = cross_val_score(model, x_test, y_test, cv=5, scor
    print("Cross validation Accuracy score: {}{}".format(round(cv_accura

    cv_precision_score = cross_val_score(model, x_test, y_test, cv=5, sc
    print("Cross validation Precision score: {}{}".format(round(cv_preci

    cv_recall_score = cross_val_score(model, x_test, y_test, cv=5, scori
    print("Cross validation Recall score: {}{}".format(round(cv_recall_s

    cv_f1_score = cross_val_score(model, x_test, y_test, cv=5, scoring="
    print("Cross validation F1 score: {}{}".format(round(cv_f1_score*100

    cf_matrix = confusion_matrix(y_test, y_pred)

    fig, ax = plt.subplots(figsize = (10,8))
    ax = sns.heatmap(cf_matrix/np.sum(cf_matrix), annot=True, fmt='.2%',

    ax.set_title('Confusion Matrix with labels\n\n', loc='left')
    ax.set_xlabel('\nPredicted Values')
    ax.set_ylabel('Actual Values ')
    ax.xaxis.set_ticklabels(['False','True'])
    ax.yaxis.set_ticklabels(['False','True'])
    plt.show()

    return model
```

## Model 1 Random Forest Classification

In [13]:

```python
# split the data
x_train, x_test, y_train, y_test = train_test_split(features, target,
                                            test_size =0.30, random_
```

In [14]: ▶
```python
1 rfm = RandomForestClassifier()
2 model1 = results(rfm, x_train, y_train, x_test, y_test)
```
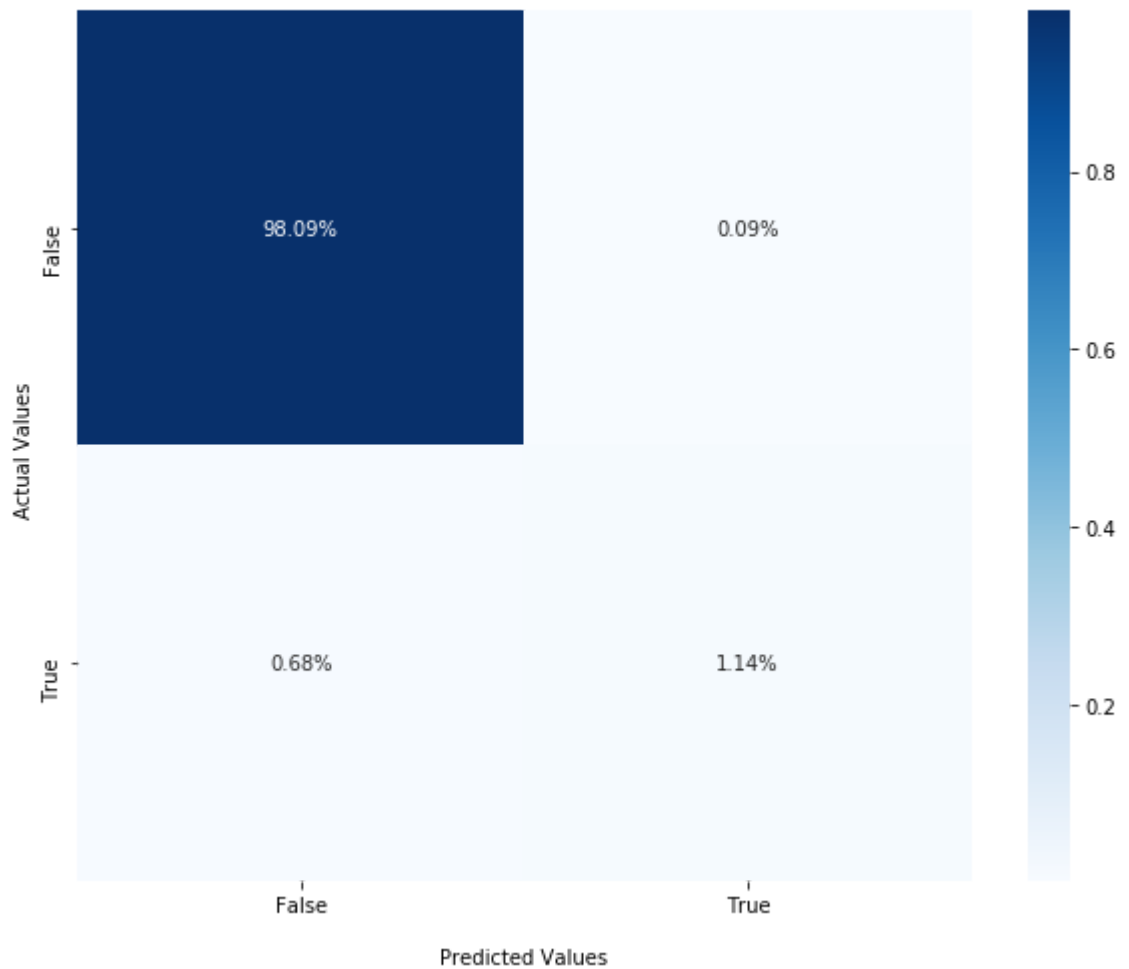
Accuracy score: 99.23%
Cross validation Accuracy score: 98.81%
Cross validation Precision score: 92.69%
Cross validation Recall score: 37.87%
Cross validation F1 score: 53.74%

Confusion Matrix with labels



**Model 2 Support Vector Machines: RBF Kernel**

In [17]: ▶
```python
1 one_hot_encoded_data = pd.get_dummies(features, columns = ['DayNight', '
2                                                            'WAT_LAND', '
```

In [18]: ▶
```python
1 # scaling the features
2 sc = StandardScaler()
3 x_train_std = sc.fit_transform(x_train)
4 x_test_std = sc.transform(x_test)
```

```
In [19]:   ▶  1  rbf = svm.SVC(kernel='rbf', C=1.0)
               2  model2 = results(rbf, x_train_std, y_train, x_test_std, y_test)
```
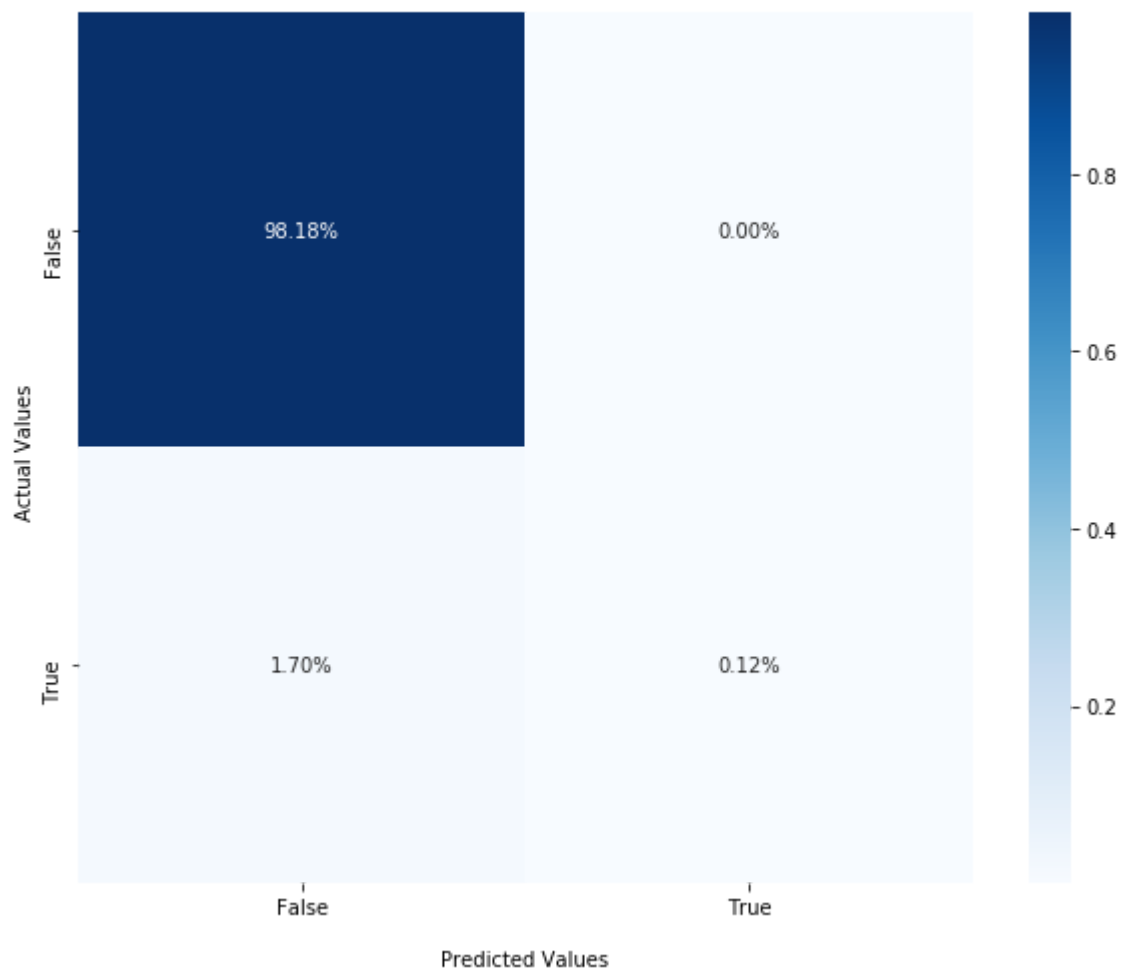
```
Accuracy score: 98.3%
Cross validation Accuracy score: 98.26%
Cross validation Precision score: 86.29%
Cross validation Recall score: 5.43%
Cross validation F1 score: 10.18%
```

Confusion Matrix with labels

# Appendix A.6 Data Modeling Preliminary LSTM RNN - Time Series Classification

In [3]:

```python
import datetime as dt
from pathlib import Path
import math
import os
import json

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.model_selection import train_test_split


from sklearn import svm
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, classification_report, confu

from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2, f_classif, mutual_info_class
from functools import partial

from sklearn.model_selection import cross_val_score
from sklearn.model_selection import RepeatedKFold
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_
from scipy import stats
from sklearn.preprocessing import StandardScaler

# Import Keras
from keras.models import Sequential
from keras.layers import Dense, LSTM, BatchNormalization
from keras.optimizers import Adam
from keras.callbacks import ReduceLROnPlateau, EarlyStopping
from keras.regularizers import l2
from time import time
import pickle

import warnings
warnings.filterwarnings('ignore')
```

In [2]: ▶|

```python
1  # Checking data type
2  def Datatype(df):
3      # shape and data types of the data
4      print("There are {} rows and {} columns".format(df.shape[0], df.shape
5      print(df.dtypes)
6
7      # select numeric columns
8      df_numeric = df.select_dtypes(include=[np.number])
9      numeric_cols = df_numeric.columns.values
10     print(numeric_cols)
11
12     # select non numeric columns
13     df_non_numeric = df.select_dtypes(exclude=[np.number])
14     non_numeric_cols = df_non_numeric.columns.values
15     print(non_numeric_cols)
```

**Functions for Feature Engineering**

In [18]: ▶|

```python
1  def transformed_features(x, y, T):
2      # scaling the features
3      sc = StandardScaler()
4      x_std = pd.DataFrame(sc.fit_transform(x.values),
5                                  index=x.index,
6                                  columns=x.columns)
7      x_std, y_std= TimeSeries_data(x_std, y, T)
8      return x_std, y_std
```

In [19]: ▶|

```python
1  def TimeSeries_data (x, y, T):
2      inputs, target = [], []
3      for i in range(y.shape[0] - (T)):
4          inputs.append(x.iloc[i:i+T].values)
5          target.append(y.iloc[i + (T)])
6      inputs, target = np.array(inputs), np.array(target).reshape(-1,1)
7      return inputs, target
```

```python
In [20]:  ▶|
     1  def lstm_transformed_data(features, target, Year_column, T):
     2      cor_matrix = features.corr().abs()
     3      upper_tri = cor_matrix.where(np.triu(np.ones(cor_matrix.shape),k=1).
     4      to_drop = [column for column in upper_tri.columns if any(upper_tri[co
     5      print("Dropping these variables because of multicollinearity")
     6      print("\n")
     7      print(to_drop)
     8
     9      features = features.drop(to_drop, axis=1)
    10      features['Year'] = Year_column['ActiveYear']
    11      features['Target'] = target['Target']
    12
    13      # split time wise because of time series classlification
    14      train_set = features[features['Year'] <2017]
    15      validation_set = features[(features['Year'] >=2017) & (features['Yea
    16      test_set = features[features['Year'] >2018]
    17
    18      x_train = train_set.drop(['Target', 'Year'], axis = 1)
    19      y_train = train_set['Target']
    20
    21      x_val = validation_set.drop(['Target', 'Year'], axis = 1)
    22      y_val = validation_set['Target']
    23
    24      x_test = test_set.drop(['Target', 'Year'], axis = 1)
    25      y_test = test_set['Target']
    26      s1 = x_val.shape
    27      s2 = x_test.shape
    28      print("\nFeatures Before Prepending {} days of data:".format(T))
    29      print("Validation set: {}".format(s1))
    30      print("Test set: {}".format(s2))
    31
    32      prepend_features1 = x_train.iloc[-(T):]
    33      prepend_features2 = x_val.iloc[-(T):]
    34
    35      x_val = pd.concat([prepend_features1, x_val], axis=0)
    36      x_test = pd.concat([prepend_features2, x_test], axis=0)
    37      s3 = x_val.shape
    38      s4 = x_test.shape
    39      print("\nFeatures After Prepending {} days of data:".format(T))
    40      print("Validation set: {}".format(s3))
    41      print("Test set: {}".format(s4))
    42
    43      x_train, y_train = transformed_features(x_train, y_train, T)
    44      x_val, y_val = transformed_features(x_val, y_val, T)
    45      x_test, y_test = transformed_features(x_test, y_test, T)
    46      s5 = x_train.shape
    47      s6 = x_val.shape
    48      s7 = x_test.shape
    49
    50      print("\nFeatures After Transforming and scaling the data:")
    51      print("train set: {}".format(s5))
    52      print("Validation set: {}".format(s6))
    53      print("Test set: {}".format(s7))
    54
    55      return x_train, y_train, x_val, y_val, x_test, y_test
    56
```

**Functions for Modeling**

In [21]: ▶|

```python
# Build the Model
def lstm_model(x_train, y_train, x_val, y_val, x_test, y_test, T,N, epoch
    model = Sequential()
    model.add(LSTM(input_shape=(T, N), units=4, activation='tanh',
                   recurrent_activation='hard_sigmoid', kernel_regulariz
                   recurrent_regularizer=l2(.01),
                   return_sequences=True, return_state=False))

    model.add(BatchNormalization())
    model.add(LSTM(units=4, activation='tanh', recurrent_activation='har
                   kernel_regularizer=l2(.01), recurrent_regularizer=l2(
                   return_sequences=True, return_state=False))

    model.add(BatchNormalization())
    model.add(LSTM(units=4, activation='tanh', recurrent_activation='har
                   kernel_regularizer=l2(0.1), recurrent_regularizer=l2(
                   return_sequences=False, return_state=False))

    model.add(BatchNormalization())
    model.add(Dense(units=1, activation='sigmoid'))
    # Compile the model with Adam optimizer
    model.compile(loss='binary_crossentropy',
                  metrics=['accuracy'],
                  optimizer=Adam(lr=.001))
    print(model.summary())
    lstm = results2(model, x_train, y_train, x_val, y_val, x_test, y_tes
    return lstm
```

```python
In [22]:    1  def results2(model, x_train, y_train, x_val, y_val, x_test, y_test, epoch
            2
            3       # Define a learning rate decay method:
            4      lr_decay = ReduceLROnPlateau(monitor='loss',
            5                                   patience=1, verbose=0,
            6                                   factor=0.5, min_lr=1e-8)
            7      # Define Early Stopping:
            8      early_stop = EarlyStopping(monitor='val_acc', min_delta=0,
            9                                 patience=30, verbose=1, mode='auto',
           10                                 baseline=0, restore_best_weights=True)
           11
           12      start = time()
           13      History = model.fit(x_train, y_train,
           14                          epochs=epoch,
           15                          batch_size=size,
           16                          validation_data=(x_val, y_val),
           17                          shuffle=True,verbose=0,
           18                          callbacks=[lr_decay, early_stop])
           19      print('-'*65)
           20      print(f'Training was completed in {time() - start:.2f} secs')
           21      print('-'*65)
           22
           23      print('\n')
           24      print('Score for Model Testing')
           25      print(model.evaluate(x_test,y_test))
           26
           27
           28      history_dict = History.history
           29      acc = history_dict['accuracy']
           30      val_acc = history_dict['val_accuracy']
           31      loss_values = history_dict['loss']
           32      val_loss_values = history_dict['val_loss']
           33      epochs = range(1,len(acc) + 1)
           34
           35      # Plotting metrics
           36      plt.plot(epochs, acc,  'bo', label = 'Training accuracy')
           37      plt.plot(epochs, val_acc, 'b', label = 'Validation accuracy')
           38      plt.title('Training and Validation Accuracy')
           39      plt.xlabel("Epochs")
           40      plt.ylabel("Accuracy")
           41      plt.legend()
           42      plt.figure()
           43      plt.plot(epochs, loss_values,  'bo', label = 'Training Loss')
           44      plt.plot(epochs, val_loss_values, 'b', label = 'Validation loss')
           45      plt.title('Training and Validation Loss')
           46      plt.xlabel("Epochs")
           47      plt.ylabel("Loss")
           48      plt.legend()
           49      plt.show()
           50
           51
           52      y_preds = model.predict_classes(x_test)
           53      accuracy = accuracy_score(y_test, y_preds)*100
           54      print("Accuracy score: {}{}".format(round(accuracy, 2), '%'))
           55
           56      precision = precision_score(y_test, y_preds)*100
```

```
57        print("Precision score: {}{}".format(round(precision, 2), "%"))
58
59        recall = recall_score(y_test, y_preds)*100
60        print("Recall score: {}{}".format(round(recall, 2), "%"))
61
62
63        cf_matrix = confusion_matrix(y_test, y_preds)
64
65        fig, ax = plt.subplots(figsize = (10,8))
66        ax = sns.heatmap(cf_matrix/np.sum(cf_matrix), annot=True, fmt='.2%',
67
68        ax.set_title('Confusion Matrix with labels\n\n', loc='left')
69        ax.set_xlabel('\nPredicted Values')
70        ax.set_ylabel('Actual Values ')
71        ax.xaxis.set_ticklabels(['False','True'])
72        ax.yaxis.set_ticklabels(['False','True'])
73        plt.show()
74
75        return model
76
```

In [23]:
```
1  def get_duration(df):
2      df['Active_minus_FireDate'] = (df["ActiveDate"] - df["FireDate"]).dt
3      df['Area_diff'] = (df["TotalAcres_sq_km"] - df["fire_dist"])
4      df = df.sort_values('Active_minus_FireDate', ascending=True)
5
6      return df
```

In [106]:
```
1  features2 = new_df.drop(['ActiveYear', 'ActiveDay','Station_dist', 'lat'
2                           'elevation','ConfidenceBinned', 'Drought_dist'
3
4  print(features.shape)
5  Year_column2 = new_df[['ActiveYear']]
6  target2 = new_df[['Target']]
7  print(target.shape)
```

```
(103121, 41)
(103121, 1)
```

**Model training for 3 day Sequential Data**

In [99]: ▶| 1  x_train, y_train, x_val, y_val, x_test, y_test = lstm_transformed_data(f

Dropping these variables because of multicollinearity


['NasaLongitude', 'Track', 'Avg_Temp', 'PS', 'TS', 'slope2', 'slope6', 'asp
ectE', 'aspectS', 'aspectW', 'CULTIR_LAND', 'SQ2', 'SQ3', 'SQ4']

Features Before Prepending 3 days of data:
Validation set: (40048, 27)
Test set: (8879, 27)

Features After Prepending 3 days of data:
Validation set: (40051, 27)
Test set: (8882, 27)

Features After Transforming and scaling the data:
train set: (54191, 3, 27)
Validation set: (40045, 3, 27)
Test set: (8876, 3, 27)

In [102]: ▶|   1   `lstm1 = lstm_model(x_train, y_train, x_val, y_val, x_test, y_test,3,27,`

Model: "sequential_3"

```
_____
Layer (type)                 Output Shape              Param #
=================================================================
lstm_7 (LSTM)                (None, 3, 4)              512
_____
batch_normalization_7 (Batch (None, 3, 4)              16
_____
lstm_8 (LSTM)                (None, 3, 4)              144
_____
batch_normalization_8 (Batch (None, 3, 4)              16
_____
lstm_9 (LSTM)                (None, 4)                 144
_____
batch_normalization_9 (Batch (None, 4)                 16
_____
dense_3 (Dense)              (None, 1)                 5
=================================================================
Total params: 853
Trainable params: 829
Non-trainable params: 24
_____
None
-----------------------------------------------------------------
Training was completed in 332.24 secs
-----------------------------------------------------------------


Score for Model Testing
8876/8876 [==============================] - 0s 43us/step
[0.845916429253353, 0.6436457633972168]
```


Training and Validation Accuracy

Training and Validation Loss

Accuracy score: 64.36%
Precision score: 22.19%
Recall score: 16.74%

Confusion Matrix with labels



**Model training for 5 day Sequential Data**

In [112]: ▶|   1  x_train, y_train, x_val, y_val, x_test, y_test = lstm_transformed_data(f

Dropping these variables because of multicollinearity

['NasaLongitude', 'Track', 'Avg_Temp', 'PS', 'TS', 'slope2', 'slope6', 'asp
ectE', 'aspectS', 'aspectW', 'CULTIR_LAND', 'SQ2', 'SQ3', 'SQ4']

Features Before Prepending 5 days of data:
Validation set: (40048, 27)
Test set: (8879, 27)

Features After Prepending 5 days of data:
Validation set: (40053, 27)
Test set: (8884, 27)

Features After Transforming and scaling the data:
train set: (54189, 5, 27)
Validation set: (40043, 5, 27)
Test set: (8874, 5, 27)

In [113]: ▶|    1   `lstm2 = lstm_model(x_train, y_train, x_val, y_val, x_test, y_test,5,27,`

```
Model: "sequential_5"
_____
Layer (type)                 Output Shape              Param #
=================================================================
lstm_13 (LSTM)               (None, 5, 4)              512
_____
batch_normalization_13 (Batc (None, 5, 4)              16
_____
lstm_14 (LSTM)               (None, 5, 4)              144
_____
batch_normalization_14 (Batc (None, 5, 4)              16
_____
lstm_15 (LSTM)               (None, 4)                 144
_____
batch_normalization_15 (Batc (None, 4)                 16
_____
dense_5 (Dense)              (None, 1)                 5
=================================================================
Total params: 853
Trainable params: 829
Non-trainable params: 24
_____
None
-----------------------------------------------------------------
Training was completed in 497.81 secs
-----------------------------------------------------------------


Score for Model Testing
8874/8874 [==============================] - 0s 55us/step
[0.6974002486592293, 0.6708361506462097]
```



Training and Validation Accuracy

Training and Validation Loss

Accuracy score: 67.08%
Precision score: 28.86%
Recall score: 21.23%

Confusion Matrix with labels



**Model training for 7 day Sequential Data**

In [114]: ▶| `1  x_train, y_train, x_val, y_val, x_test, y_test = lstm_transformed_data(f`

Dropping these variables because of multicollinearity

['NasaLongitude', 'Track', 'Avg_Temp', 'PS', 'TS', 'slope2', 'slope6', 'asp
ectE', 'aspectS', 'aspectW', 'CULTIR_LAND', 'SQ2', 'SQ3', 'SQ4']

Features Before Prepending 7 days of data:
Validation set: (40048, 27)
Test set: (8879, 27)

Features After Prepending 7 days of data:
Validation set: (40055, 27)
Test set: (8886, 27)

Features After Transforming and scaling the data:
train set: (54187, 7, 27)
Validation set: (40041, 7, 27)
Test set: (8872, 7, 27)

In [114]: ▶| `1  x_train, y_train, x_val, y_val, x_test, y_test = lstm_transformed_data(f`

In [115]:  ▶|    1  lstm3 = lstm_model(x_train, y_train, x_val, y_val, x_test, y_test,7,27,

```
Model: "sequential_6"
_____
Layer (type)                 Output Shape              Param #
=================================================================
lstm_16 (LSTM)               (None, 7, 4)              512
_____
batch_normalization_16 (Batc (None, 7, 4)              16
_____
lstm_17 (LSTM)               (None, 7, 4)              144
_____
batch_normalization_17 (Batc (None, 7, 4)              16
_____
lstm_18 (LSTM)               (None, 4)                 144
_____
batch_normalization_18 (Batc (None, 4)                 16
_____
dense_6 (Dense)              (None, 1)                 5
=================================================================
Total params: 853
Trainable params: 829
Non-trainable params: 24
_____
None
-----------------------------------------------------------------
Training was completed in 724.57 secs
-----------------------------------------------------------------


Score for Model Testing
8872/8872 [==============================] - 1s 73us/step
[1.2699199670655326, 0.6459648609161377]
```



Training and Validation Accuracy

## Training and Validation Loss



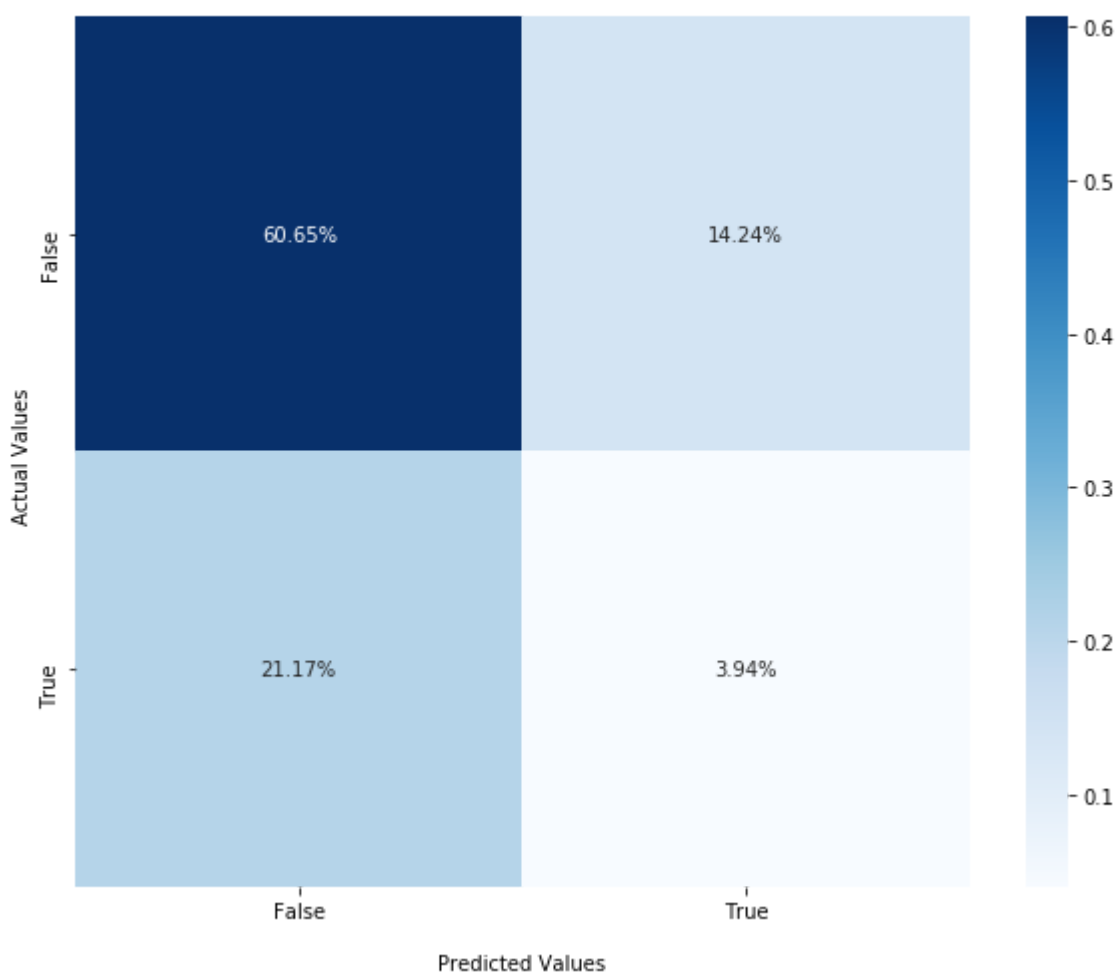Accuracy score: 64.6%
Precision score: 21.7%
Recall score: 15.71%

## Confusion Matrix with labels



In [ ]:  ▶|    1