

# Analyzing California Wildfire

October 2, 2022

## Goal

Rising temperatures, drought conditions, and dry vegetation are all contributing causes of wildfires in California, but the worsening impact of climate change is further exacerbating and prolonging the fire season in the state.[1]

As part of this assignment, I analyze the wildfire in California to understand the trends of wildfire from 2013-2020.[2] Using a open-dataset from Kaggle, I have analyzed the distributions of wildfires, the impact wildfires have and the amount of resources that are used during a wildfire.

Github: <https://github.com/maitreyi-kv/cmpe255>

## Dataset Details

The collection includes information about the counties where wildfires have occurred, their latitude and longitude values, and the times the flames started. With this dataset the goal is to obtain insights on the areas of California that are most at risk for wildfires, as well as when and how frequently they occur.

Dataset Link: <https://www.kaggle.com/datasets/ananthu017/california-wildfire-incidents-20132020>

This dataset is uploaded into big query in dataset *california* under the table *cali\_fire*. The result dataframe after processing is stored in table *table\_fire\_post\_processing*.

Row	County	Lat	Lon	CAL_FIRE_Series_County	ACTIVITY	ACTIVITYTYPE	YEAR	INCIDENTID	INCIDENTID_Series_Fire	Coordinates_Latitude	Coordinates_Longitude
1	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611
2	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611
3	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611
4	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611
5	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611
6	San Diego	32.7157	-117.1611	CAL_FIRE_Series_County_1	fire	Incident	2013	Incidents20131218-1000000000	Incidents20131218-1000000000	32.7157	-117.1611

Figure 1: Dataset before and post processing on Big Query

# Loading Data, Analysis and Observations

## 0.1 Loading Data

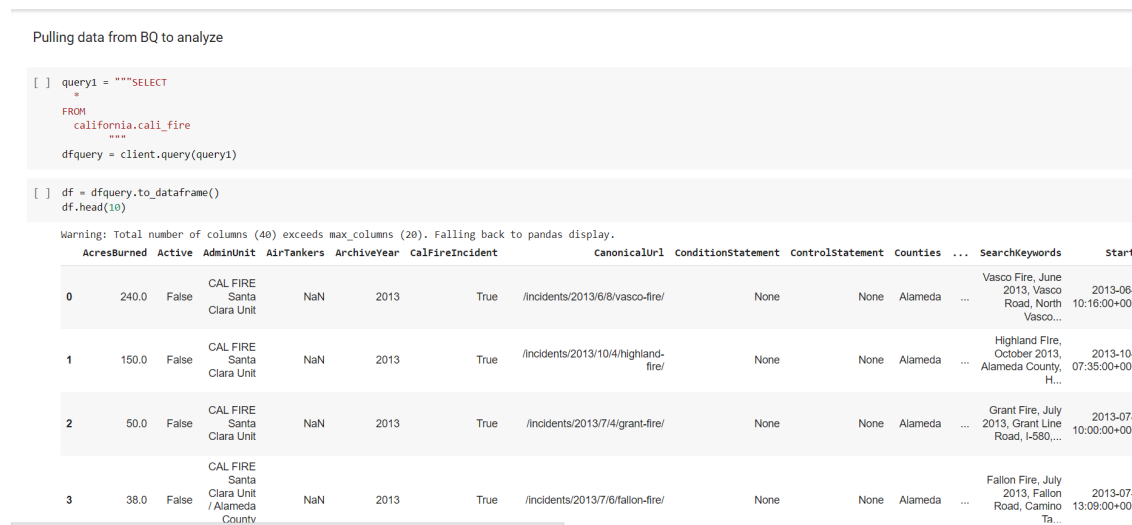


Figure 2: Dataset loaded in Colab from BQ

## 0.2 Basic Analysis

The dataset has 1636 wildfire incidents over the course of 2013 to 2019. There are 40 features for the incident. All the wildfires in the dataset are inactive.

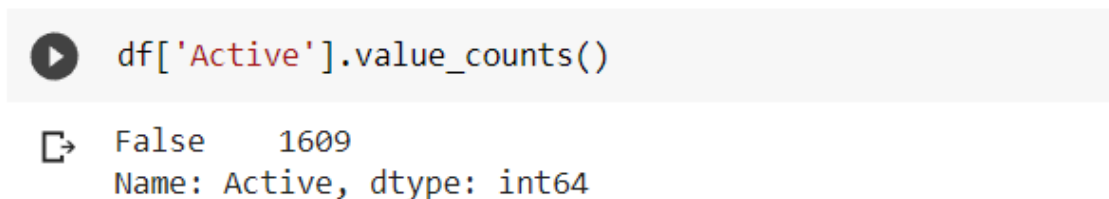


Figure 3: All wildfires are inactive

### 0.3 Preprocessing

I am performing preprocessing to clean the dataset. Broadly they belong to the following buckets.

1. Duration - In some of the records, the end-time is lesser than the start time. In some other cases, the end-time is so big indicating a super long wild-fire (verified that these are inaccurate). I am dropping these values as it is a mistake in the data entry.

```
df[df['Duration'].describe()

count    1550.000000
mean      84.549758
std       883.374905
min     -17052.000000
25%        1.750000
50%       19.958333
75%      170.822917
max     17900.708333
Name: Duration, dtype: float64

df[df.Duration < 0]['Duration'].count()

28

IncorrectDuration = df[(df.Duration < 0) | (df.Duration > 200)]['Duration']

IncorrectDuration.count()

235
```

Figure 4: Negative duration - incorrect data entry

2. Lat/Long - In some of the records, I can observe that the latitudes and longitudes are invalid.

```
df[['Latitude', 'Longitude']].describe()

index      Latitude      Longitude
count      1374.0
mean      37.476997280932395      -106.70279755922448
std       147.70608091372816      38.86836444970496
min       -120.258
25%      34.166816999999995      -121.79006975
50%      37.2225835
75%      39.16737
max       5487.0      118.9082

df.drop(df[(df.Latitude < -90) | (df.Latitude > 90)].index, inplace=True)
df.drop(df[(df.Longitude < -180) | (df.Longitude >= 0)].index, inplace=True)
```

Figure 5: Negative duration - incorrect data entry

3. Duplicates - There are duplicate records in the data-set.



```
df = df.drop_duplicates(subset=['Name', 'Started', 'AcresBurned', 'StructuresDamaged', 'StructuresDestroyed'], keep='first', inplace=False, ignore_index=False).reset_index().drop(column
```

Figure 6: Removing duplicate entries

Records in the data-set which belong to either of the above buckets are dropped.

## 0.4 Observations and Visualization

1. How long do wildfires last for? Analyzing the duration of the wildfires to understand length.

As we can see most fires last less then a day. Long lasting fires are more more less frequently occurring than shorter fires.

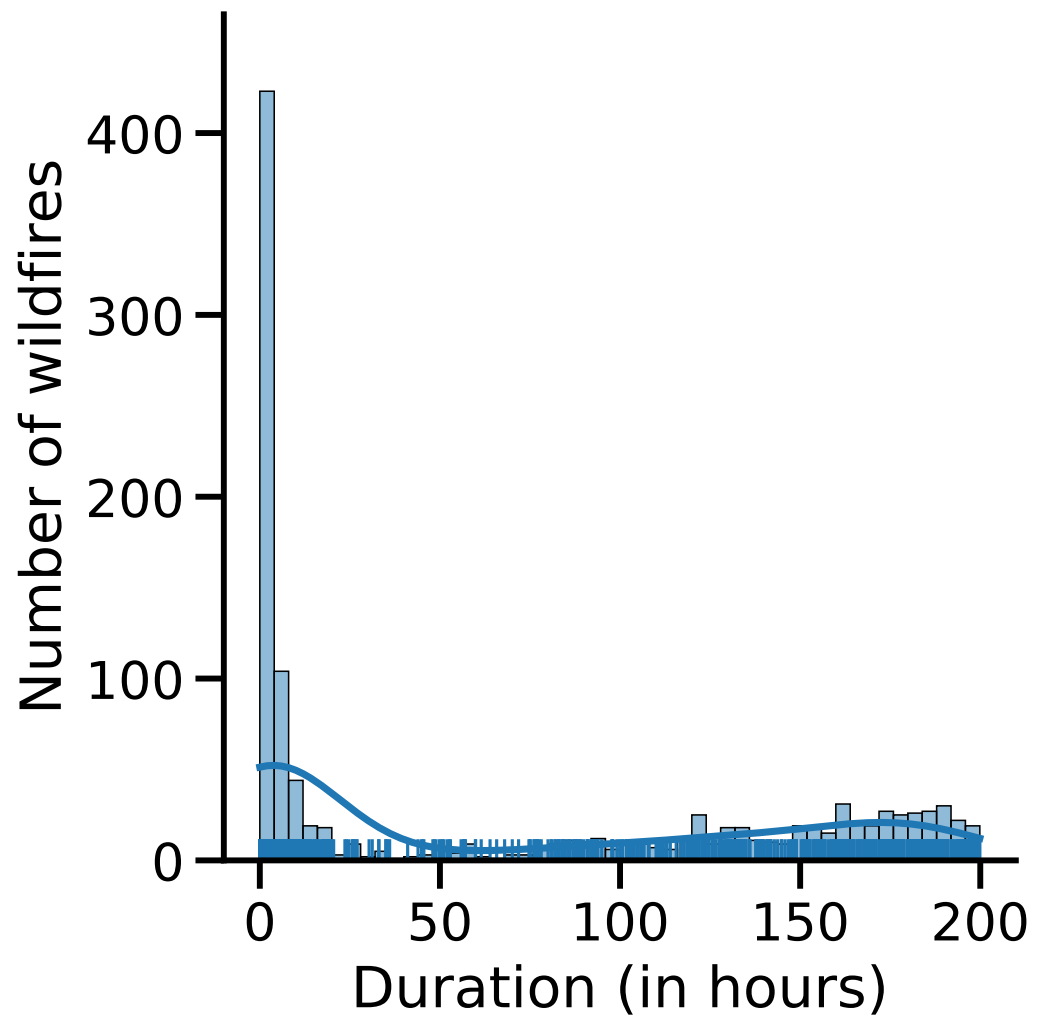


Figure 7: Length of wildfire

2. Where do wildfires occur? Analyzing the locations of the wildfires to understand if any area is more prone to wildfire.

Riverside and San Diego are places that are most affected since the last 6 years by wildfire.

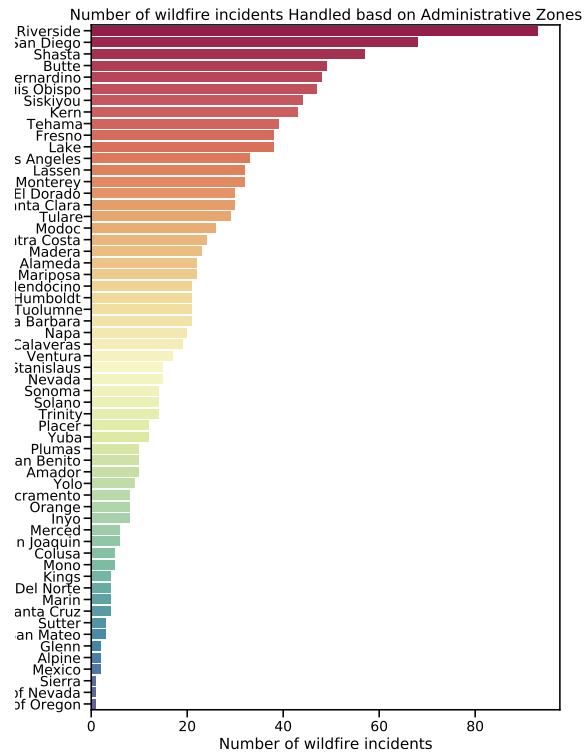


Figure 8: Most affected counties

### 3. Analyzing Administrative unit performing relief

There are 417 admin units in the dataset.

Looking at the administrative zones for the California Fire Dept. (shown below) we observe that the Sanoma Lake, Santa Clara, and Shasta-Trinity Unit are the top three administrative zones dealing with fire.

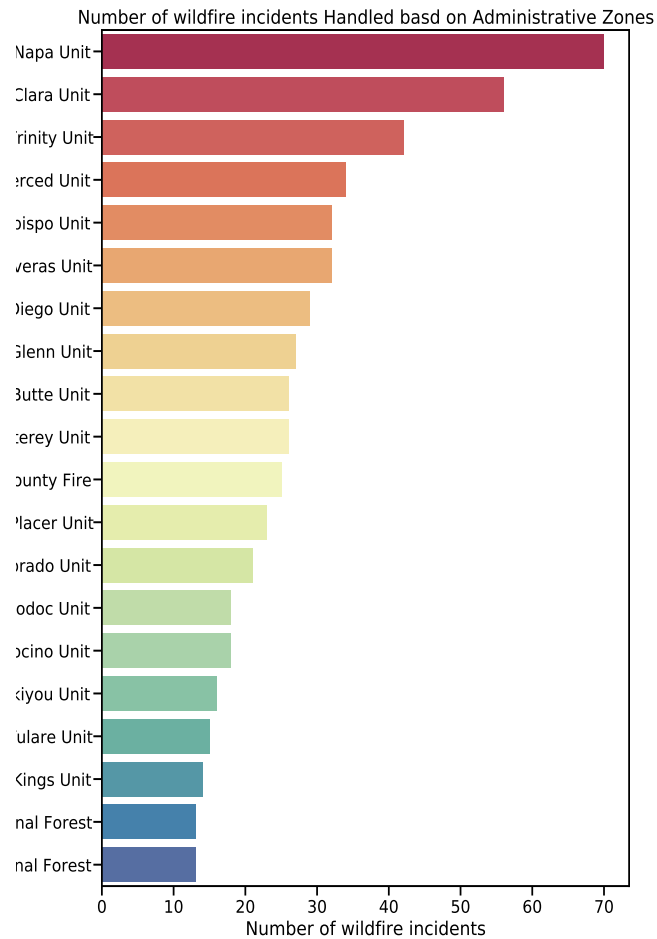


Figure 9: Top administrative unit performing relief

#### 4. Analyzing geolocation of areas burned in fire incidents

Major wildfire incidents have mostly occurred near regions of large population density and mostly concentrated towards South Eastern regions (San Diego and Riverside) and above San Francisco.

Wildfires since 2017 have been mostly concentrated in the same regions in California.

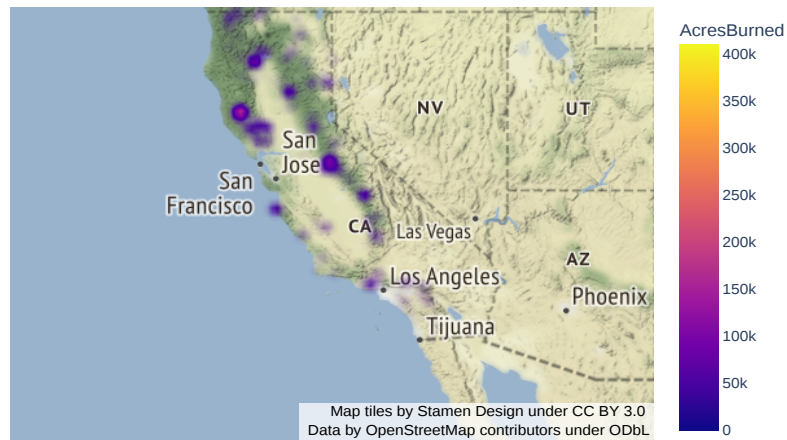


Figure 10: Density (in acres) of area burned by Wildfire



5. How many personnel's are involved in fighting Wildfires. How much fatality occurs across locations. Both personnel's and fatalities data correlates with the locations which have been heavily affected by wildfires.

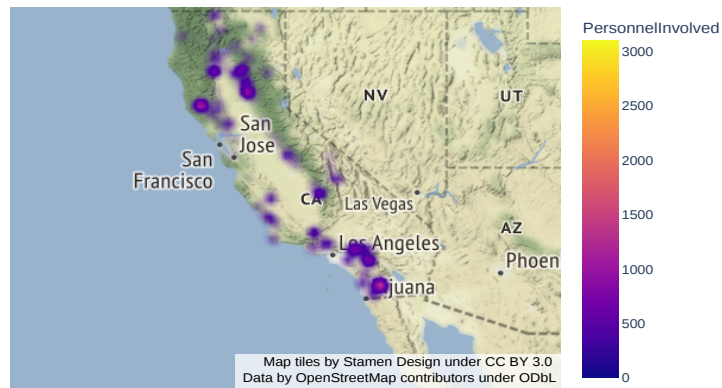


Figure 11: Personall's putting out fire in locations



Figure 12: Fatalities occurring in locations

6. How many wildfires occur each year?

2017 has seen the most number of wildfires. From 2014 there is a constant increasing trend which could be due to global warming and climate change.

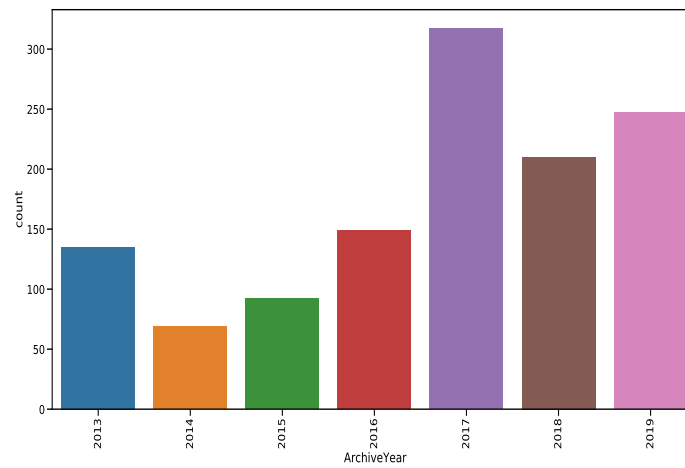


Figure 13: Wildfire by year

7. Are the acres burned directly related to no of wildfires.

Although there is a positive correlation, it is not always true as the largest wildfire is in 2018. The number of wildfires is lesser than 2017 yet the acres burned is much more.

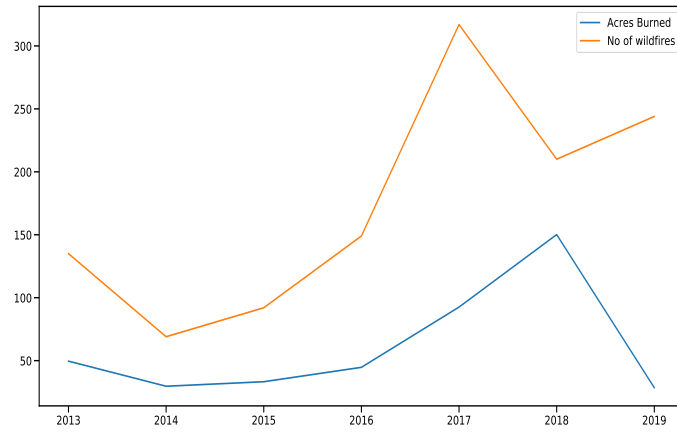


Figure 14: Correlation between wildfire and acres burned

8. What about correlation between acres burned and water tenders? And correlation between acres burned and personnel's? Are there any patterns/changes across the years?

The water tenders are is correlated. Especially during the major fire event in 2018 - most number of water tenders are used. Although, in the beginning years the number is much more. This decreases over time, and less water tender is being used. This could be due to technological advancement in mitigating wildfires. Similar patterns are seen in personnel's, they decrease over time which could be due to tech advancements.

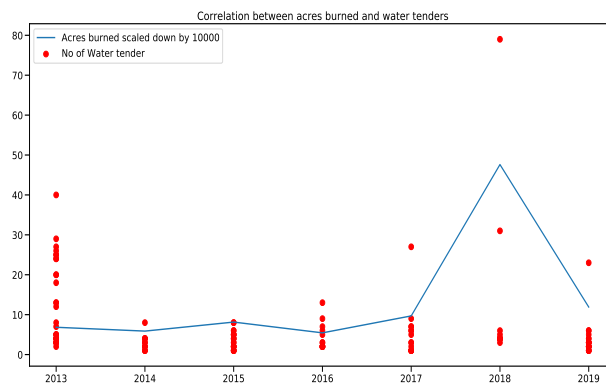


Figure 15: Correlation between wildfire and water tender

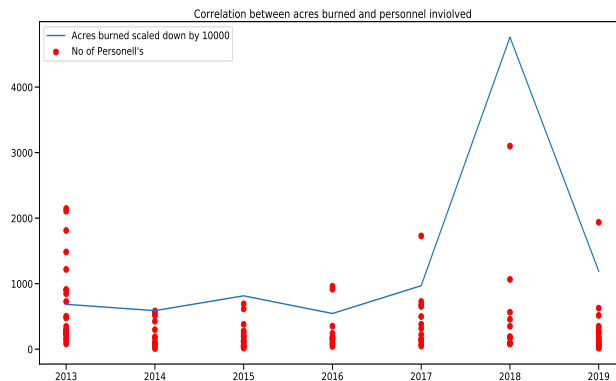


Figure 16: Correlation between wildfire and personell's

9. What is the correlation amongst all features? Is there any high correlation features in this?

The Pearson's correlation shows that there are high correlations between

1. corr of .90 between Personell's involved and Fatalities/Engines 2. corr of .99

There is Moderate correlation between Water tenders and personnel's involved. This can be visualized in the pair plot below too.

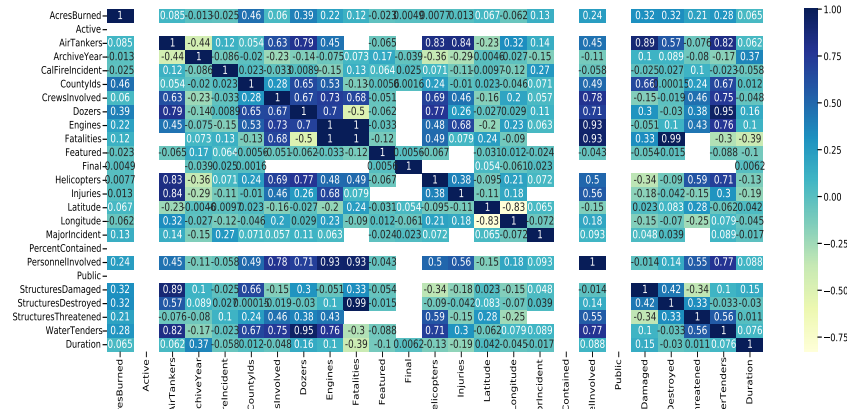


Figure 17: Pearsons correlation

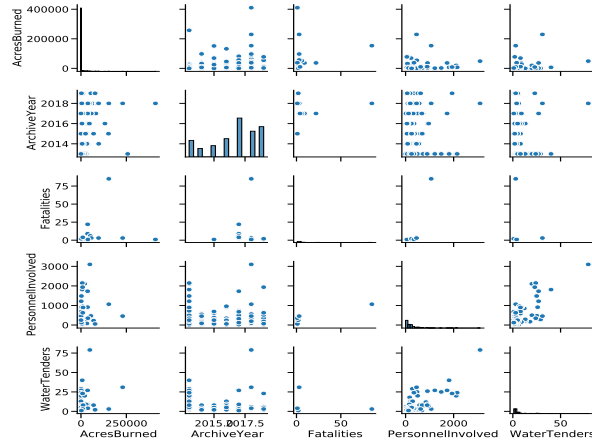


Figure 18: Pair Plots

## 0.5 Webapp

Using plotly and dash, I have created a app to visualize 4 different charts - Geo Location Map, Histogram, Scatter Plot, Pie Chart

1. Locations affected by wildfires using Geo Location Map
2. Counties impacted by wildfires using Histogram
3. Top 5 administrative units performing relief using Histogram
4. Distribution of wildfire by county using Scatter Plot
5. Distribution of wildfire by year using Pie Chart

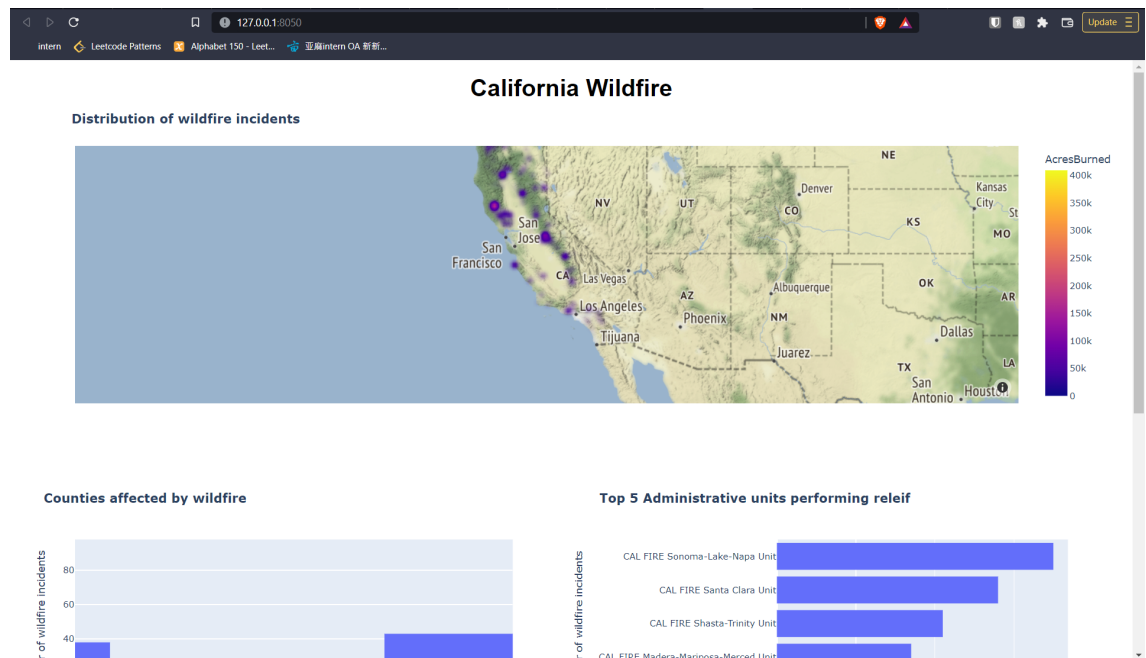


Figure 19: Webapp using plotly and Dash, Page 1

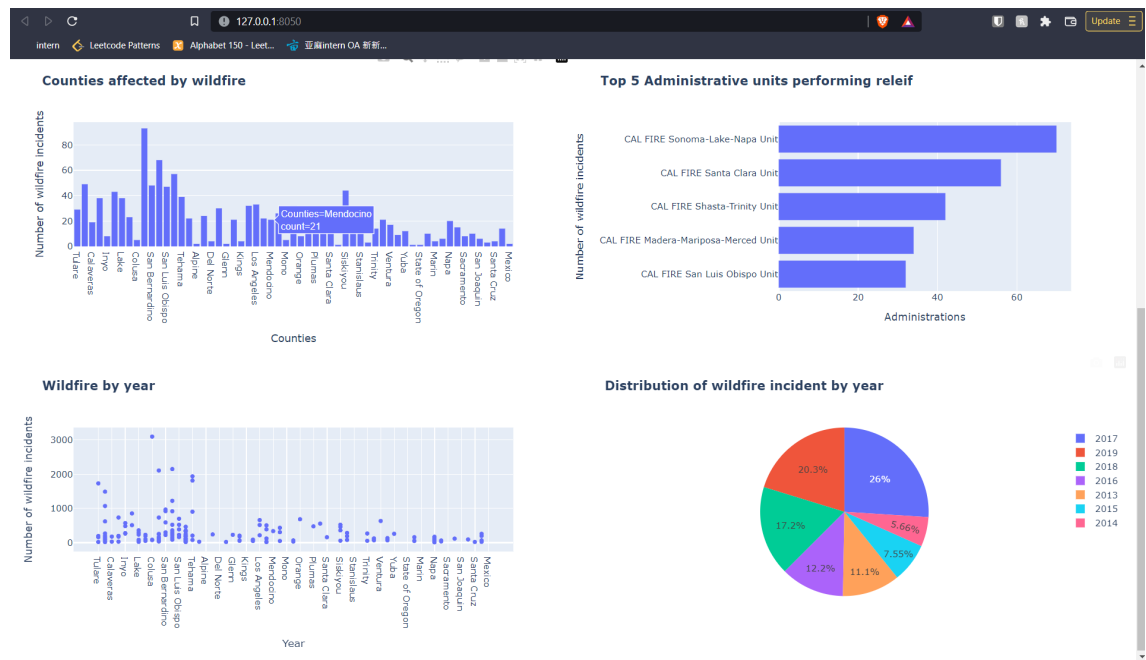


Figure 20: Webapp using plotly and Dash, Page 2

## References

- [1] <https://earth.org/what-causes-california-wildfires/#:~:text=Rising>
  - [2] <https://www.kaggle.com/datasets/ananthu017/california-wildfire-incidents-20132020>
- <https://plotly.com/python/line-charts/>