# Predicting Forest Fires in California

## Using Historical Weather and Fire Data

By: Adam Swan, Jeff Warchall, Omar Younis

Image Source: Noah Berger - Associated Press

# Problem Statement:

To help the California Department of Forestry and Fire Protection allocate resources, can we predict the likelihood of fires utilizing historical weather and wildfire data?

Is it desirable to accurately predict as many fires as possible, or:

To minimize false alarms?

# Data Sources

- Geographic and wildfire data from the State of the California
  - 4,279 fires from July 1, 2008 to December 31, 2020

- Weather data from World Weather Online
  - Historical weather records from every day in the same range
  - Aggregated to a monthly basis

# Web Scraping - Weather

- World Weather Online:
    - Includes temperature, precipitation, and wind data for:
    - Every weather station, for:
    - Every day:
    - Since July 1, 2008

- Data obtained from every county in California
    - Aggregated by month

# Web Scraping - Fires

- The available fire data includes:
  - Date of fire
  - Number of acres burned
  - Cause of fire

- Data extends back to 1952, but truncated at 2008 based on available weather data

# Preprocessing

- Preprocessing was extensive for this project:
  - Multi-indexing
  - Monthly aggregation
  - Quarterly aggregation
  - Table joins
  - Multiple simultaneous fires
  - Data imputing
  - Time series (?)

# Preprocessing

- Multi-Indexing

| | date | county | maxtempF | mintempF | avgtempF | totalSnow_cm | humid | wind | precip | sunHour | lat | long |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 2008-07 | Sierra County | 86.290323 | 46.645161 | 76.709677 | 0.0 | 32.709677 | 5.451613 | 0.000000 | 14.396774 | 39.58 | -120.52 |
| **1** | 2008-07 | Sacramento County | 97.290323 | 63.419355 | 86.516129 | 0.0 | 39.838710 | 4.741935 | 0.000000 | 13.741935 | 38.45 | -121.34 |
| **2** | 2008-07 | Santa Barbara County | 89.129032 | 59.709677 | 80.548387 | 0.0 | 41.451613 | 7.354839 | 0.000000 | 13.164516 | 34.54 | -120.04 |
| **3** | 2008-07 | Calaveras County | 96.419355 | 51.290323 | 87.032258 | 0.0 | 33.580645 | 5.387097 | 0.000000 | 14.022581 | 38.18 | -120.56 |
| **4** | 2008-07 | Ventura County | 78.612903 | 61.354839 | 73.193548 | 0.0 | 41.548387 | 5.483871 | 0.003226 | 13.551613 | 34.36 | -119.13 |

# Preprocessing

- Aggregation
  - Data collected form every day
  - Looking for monthly / quarterly data

| | date | maxtempF | mintempF | avgtempF | totalSnow_cm | sunHour | precip | humidity | windspeed | lat | long |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **105845** | 2020-12-27 | 55 | 38 | 49 | 0.0 | 8.7 | 0.0 | 97 | 7 | 35.39 | -120.45 |
| **105846** | 2020-12-28 | 48 | 39 | 47 | 0.0 | 3.6 | 0.0 | 97 | 7 | 35.39 | -120.45 |
| **105847** | 2020-12-29 | 54 | 39 | 48 | 0.0 | 8.7 | 0.0 | 97 | 7 | 35.39 | -120.45 |
| **105848** | 2020-12-30 | 55 | 38 | 47 | 0.0 | 8.7 | 0.0 | 97 | 7 | 35.39 | -120.45 |
| **105849** | 2020-12-31 | 55 | 42 | 49 | 0.0 | 9.9 | 0.0 | 97 | 7 | 35.39 | -120.45 |

# Preprocessing

- Joining simultaneous records
  - More than 1 fire in a county in a month

| | UNIT_ID | FIRE_NAME | ALARM_DATE | CONT_DATE | CAUSE | REPORT_AC | GIS_ACRES | SHAPE_Length | SHAPE_Area |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Yuba County | NELSON | 2020-06-18 00:00:00+00:00 | 2020/06/23 00:00:00+00 | 11.0 | 110.0 | 109.60250 | 4179.743142 | 7.331347e+05 |
| 1 | Yuba County | AMORUSO | 2020-06-01 00:00:00+00:00 | 2020/06/04 00:00:00+00 | 2.0 | 670.0 | 685.58502 | 12399.375391 | 4.578172e+06 |
| 2 | Yuba County | ATHENS | 2020-08-10 00:00:00+00:00 | 2020/03/01 00:00:00+00 | 14.0 | 26.0 | 27.30048 | 2119.194120 | 1.823876e+05 |
| 3 | Yuba County | FLEMING | 2020-03-31 00:00:00+00:00 | 2020/04/01 00:00:00+00 | 9.0 | 13.0 | 12.93155 | 2029.524881 | 8.667942e+04 |
| 4 | Yuba County | MELANESE | 2020-04-14 00:00:00+00:00 | 2020/04/19 00:00:00+00 | 18.0 | 10.3 | 10.31596 | 1342.742903 | 7.017912e+04 |

# Preprocessing

- Data imputation
  - Months when there is no fire
  - Cases when a fire's cause is missing

```
Name: CAUSE, dtype: int64

[ ]: x , y = final.shape
     for i in range(x):
         clear_output()
         print(f"{i} of {x-1}")
         if final.iloc[i,-1] != 0 and pd.isnull(final.iloc[i,-2]):
             final.iloc[i,-2] = mode_cause

10750 of 11148
```

# Preprocessing

- Is this a time series?
  - Ultimately, no; but:
  - Time series tools were useful in feature engineering

# Final Dataset

- The final dataset passed to the visualization department contained:
  - 10,988 records
  - 4,297 of which had fires
  - Features for:
    - Monthly weather averages
    - Quarterly weather averages
    - Quarterly cumulative precipitation
    - Acres burned
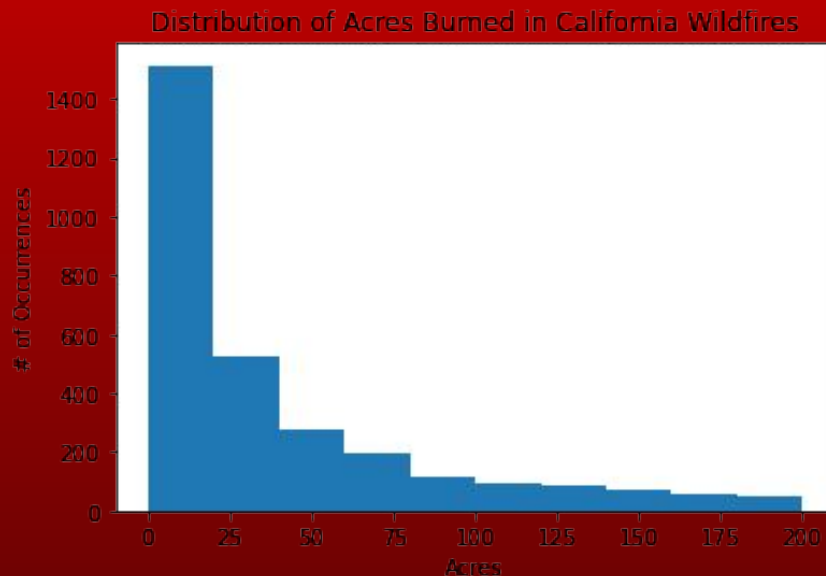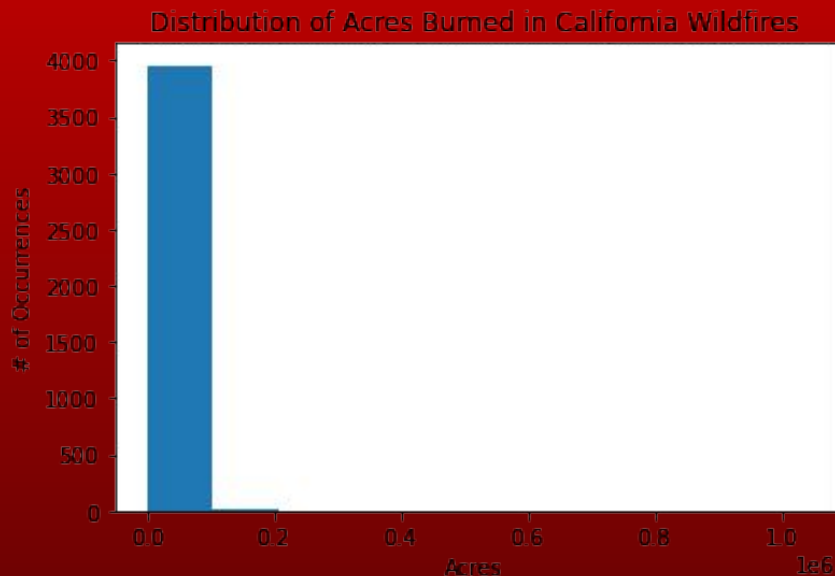    - Cause of fire
    - Geo Location

# Data Dictionary

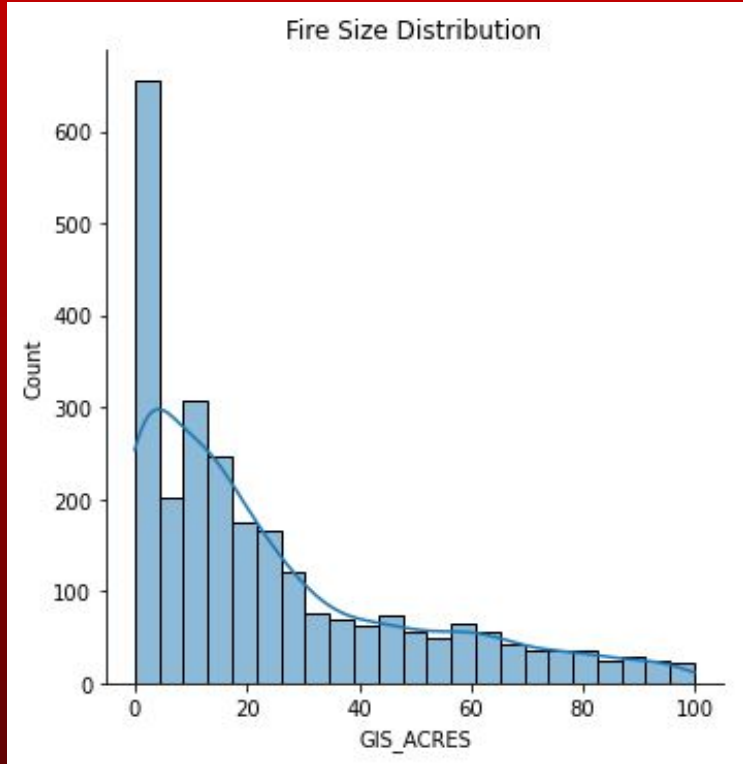| Feature | Type | Description |
|---|---|---|
| date | *object* | The month and year of when the fire took place. |
| county | *object* | The county the fire started in. |
| maxtempF | *float* | The average maximum temperature of that month in °F. |
| mintempF | *float* | The average min temperature of that month in °F. |
| avgtempF | *float* | The average average temperature of that month in °F. |
| totalSnow | *float* | The total snow for that month. |
| humid | *float* | The average humidity for that month. |
| wind | *float* | The average wind for that month. |
| precip | *float* | The average precipitation for that month. |
| q_avgtempF | *float* | The quarterly average temperature in °F. |
| q_avghumid | *float* | The quarterly average humidity. |
| q_sumprecip | *float* | The quarterly average precipitation. |
| sunHour | *float* | The average hours of sun for that month. |
| FIRE_NAME | *object* | The name of the fire. |
| CAUSE | *float* | The cause of the fire. |
| lat | *float* | The latitude coordinate of the fire's location. |
| long | *float* | The longitude coordinate of the fire's location. |
| GIS_ACRES | *float* | The total number of arces burned. |

# How Big is an Acre?


One acre is approximately 60% of a soccer pitch

One Acre

thecalculatorsite.com


One acre is approximately 16 tennis courts

One Acre

thecalculatorsite.com

43,560 sq/ft or .001562 sq/mi

Image source - www.thecalculatorsite.com

# Fire Size Distribution:



61% of 4,279 fires between 2008 - 2020 were smaller than 100 acres -  roughly .15 sq/mi

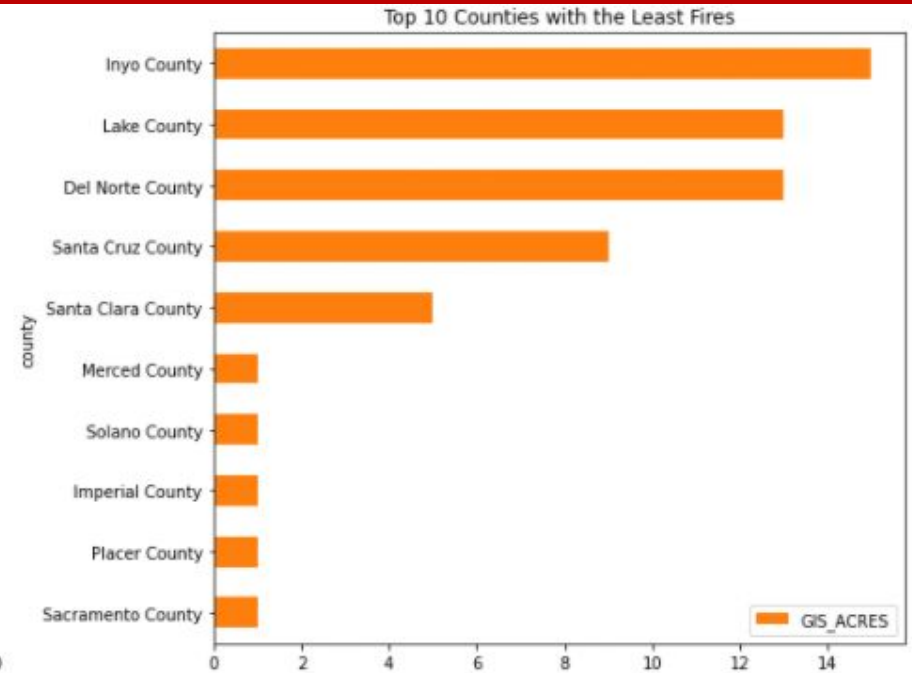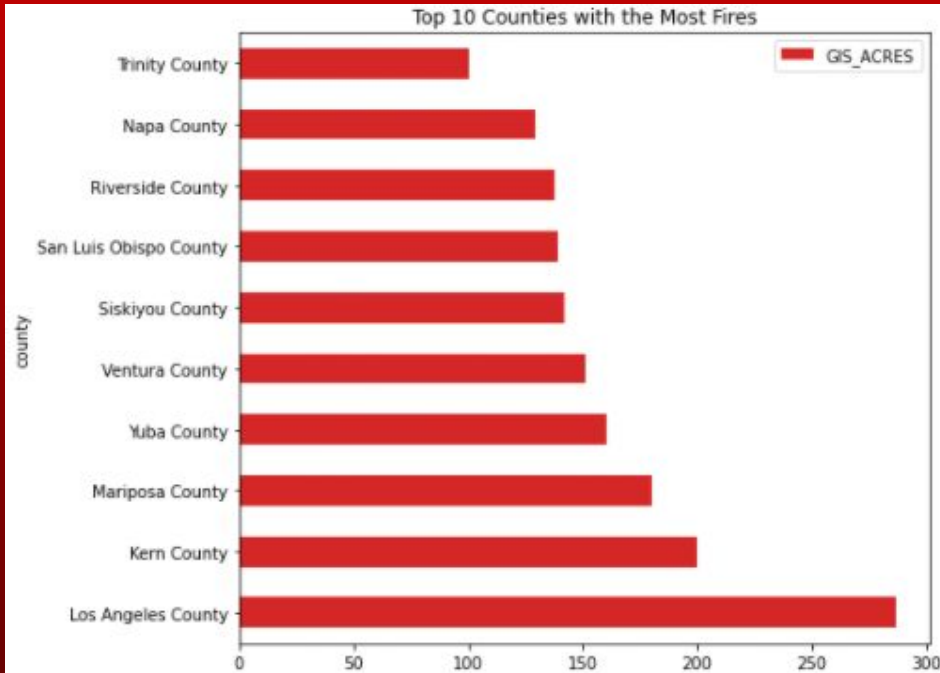# Fire Size Distribution (Continued)



- Break Down Version.

# Largest Fire:

Year - 2020
Cause - Lightning
Damage - 1.03 million
acres (1,562 sq/mi),
nearly the size of
Harris county, Texas



Harris County © Texas Almanac

# Top 8 Causes of CA Wildfires:

| Rank | Cause | # of Occurrences |
|:---:|---|---|
| **1** | Unknown | 1,208 |
| **2** | Lightning | 832 |
| **3** | Equipment use | 439 |
| **4** | Vehicle | 262 |
| **5** | Powerline | 207 |
| **6** | Arson | 184 |
| **7** | Debris | 153 |
| **8** | Campfire | 96 |

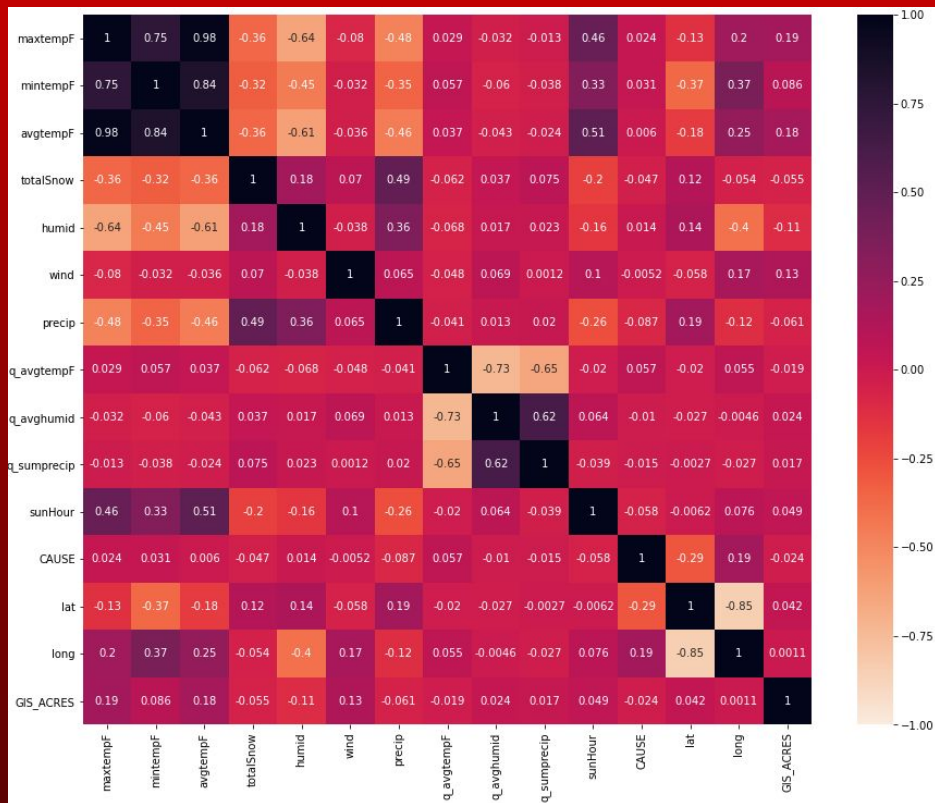# Counties with Most and Least Fires

# Interactive Fire Map



- Interactive .html file.

- Fire Size (marker size)

- All Fires (not limited to sizes less than 100 acres).

- Most fires in Northern California.
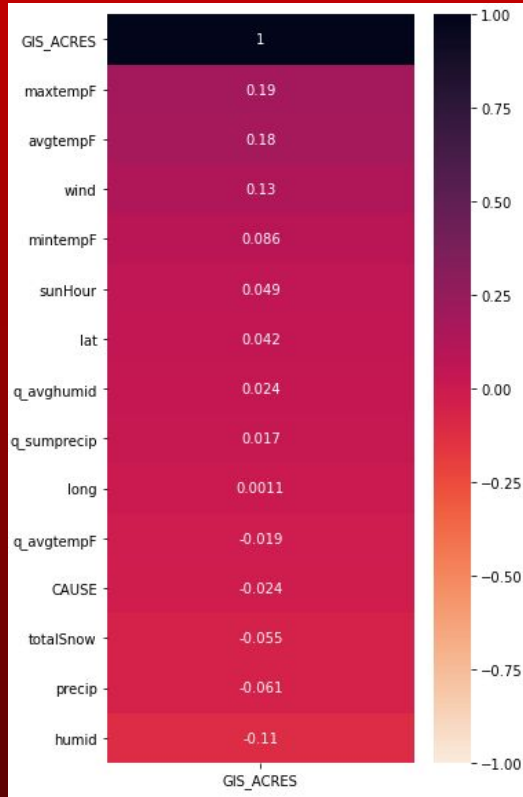
# Top 10 Under 100



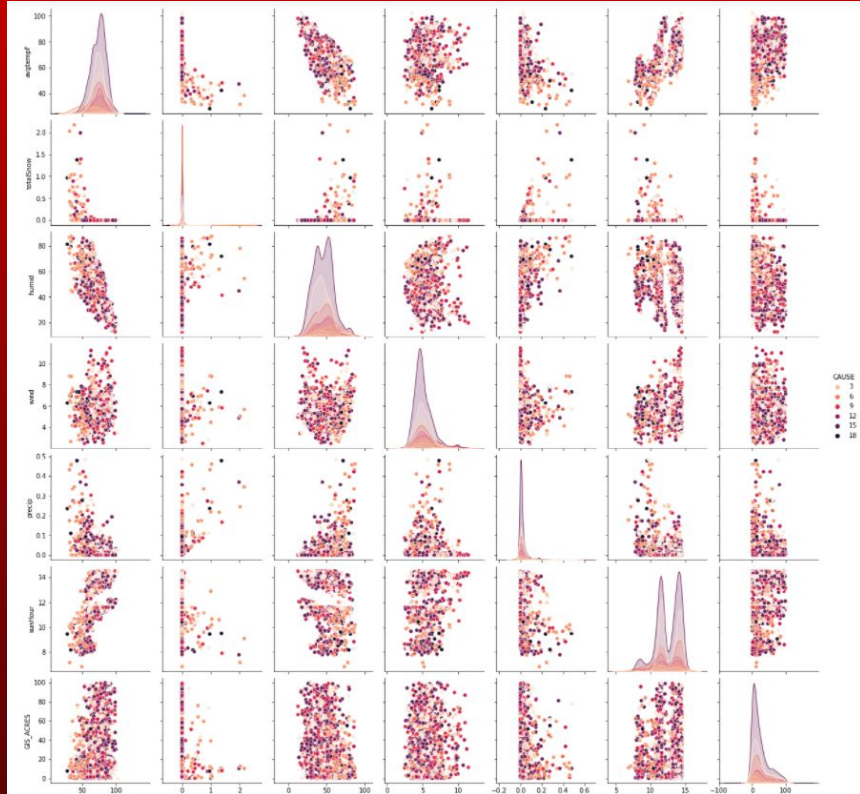| | county | FIRE_NAME | GIS_ACRES |
|---|---|---|---|
| 1827 | Riverside County | GILMAN | 99.837486 |
| 7822 | Shasta County | R-3 MUD | 99.573639 |
| 9543 | Santa Barbara County | RANGE | 99.083054 |
| 2555 | Riverside County | CABAZON | 98.980995 |
| 4373 | Napa County | HIGHLAND | 98.822746 |
| 1378 | Kern County | BRAMLETTE | 98.147751 |
| 3369 | Riverside County | FREEWAY | 97.775436 |
| 6014 | Lake County | DEER | 97.552078 |
| 2455 | Riverside County | CAHUILLA | 97.545860 |
| 6020 | Glenn County | 36 | 97.439041 |

# Heatmap (All)



- Significant correlation
  - |0.4| or greater

- No strong correlation between variables
  - Dependent Features
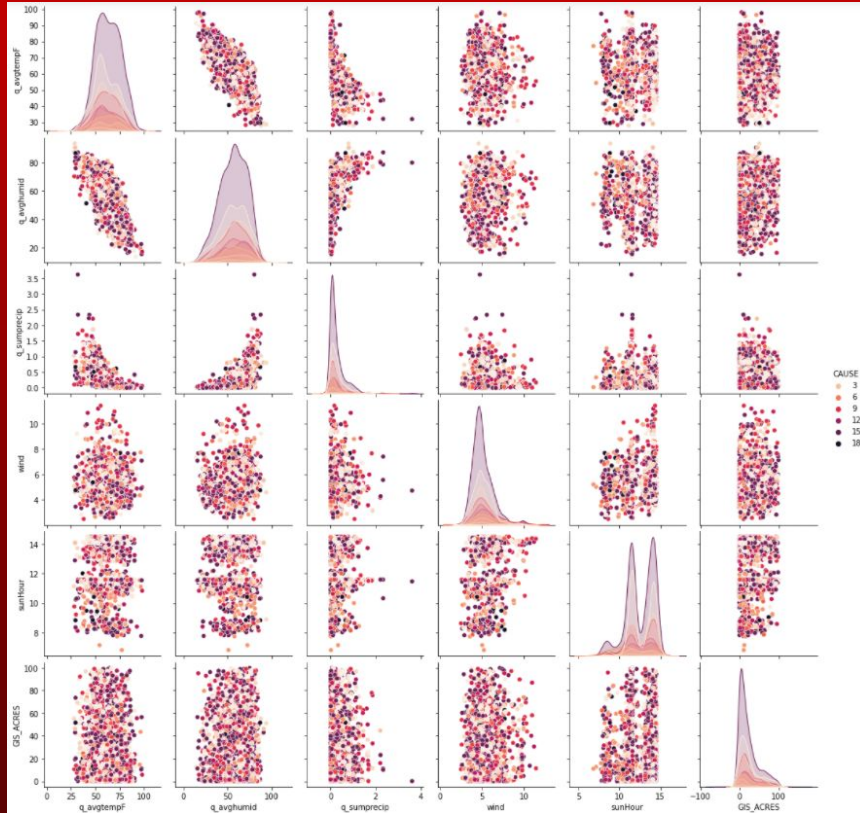
# Heatmap (GIS_ACRES)



- Significant correlation
  - |0.4| or greater

- No strong correlations

# Pairplot Monthly



- No strong patterns in the scatter plots.

- Dependant Features.

# Pairplot Monthly



- No strong patterns in the scatter plots.

- Dependant Features.

# Modeling:

1) Logistic Regression

2) KNN Classifier

3) Random Forest Classifier

4) Voting Classifier

# Model Performance:

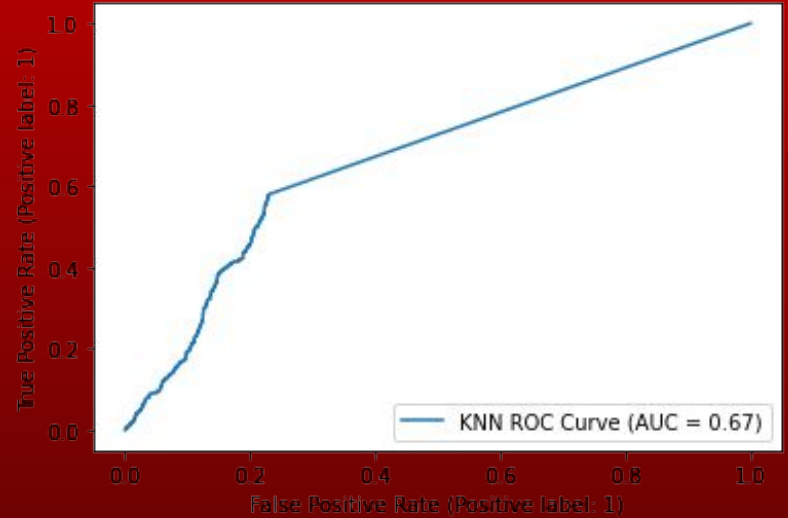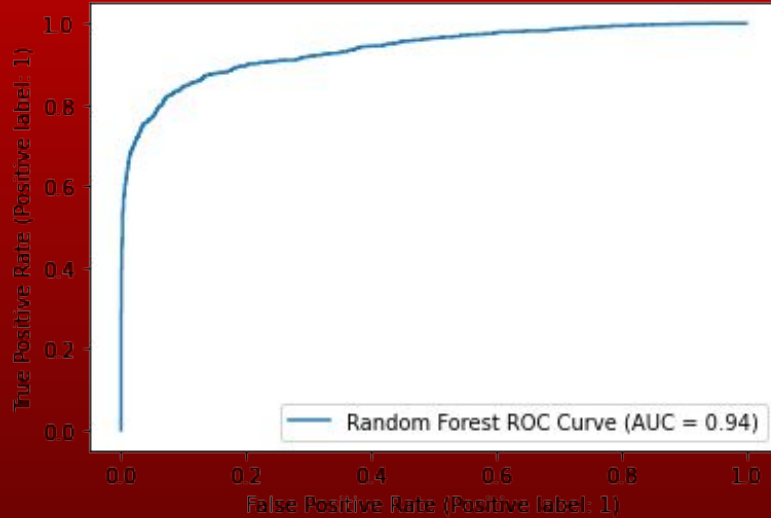| Model Type | Accuracy |
|---|---|
| Logistic Regression | 76% |
| KNN Classifier | 85% |
| → Random Forest Classifier | 88% |
| Voting Classifier | 87% |

# Logistic Regression Coefficients:

| Variable: | Effect: |
|-----------|---------|
| August | +108% |
| September | +69% |
| Hours of sun | +51% |
| November | +42% |
| July | +37% |
| December | +21% |
| Total snowfall | +2% |

| Variable: | Effect: |
|-----------|---------|
| February | -64% |
| March | -50% |
| May | -49% |
| January | -49% |
| April | -44% |
| Wind | -28% |
| June | -22% |

# Random Forest Features:

| Rank | Feature |
|------|---------|
| 1 | Hours of sunlight |
| 2 | Average temp |
| 3 | Average humidity |
| 4 | Average wind speed |
| 5 | Average precipitation |
| 6 | Year |
| 7 | July |
| 8 | Monthly snowfall |

# Additional Metrics:

# Additional Metrics cont:

*Null model: 61% accuracy*

| Random Forest Classifier | |
|---|---|
| **Metric** | **Score** |
| AUC - ROC | 87% |
| Recall | 83% |
| Precision | 84% |
| Accuracy | 88% |

| Voting Classifier | |
|---|---|
| **Metric** | **Score** |
| AUC - ROC | N/A |
| Recall | 86% |
| Precision | 80% |
| Accuracy | 87% |

**Conclusion** :

**Random Forest:**

1) Higher accuracy
2) Higher Precision

**Voting Classifier:**

1) Lower precision
2) Higher recall

## Recommendation:

### <u>Random Forest</u>

If ample resources and funding are available, we recommend the Random Forest model as it positively identifies forest fires more often than our other models.

**Cons** - More false positives, possibly wasting resources.

**Final Considerations:**

Forest fires are:
- Unavoidable
- Healthy for a forest's life cycle
- Generally small when responded to quickly
- Increasing in danger and damage as humans populate and expand further out from urban and suburban areas

# End

Thank you for your time, we will now take questions.