

# Project Proposal

Info Innovators - Kevin Mao, Arnav Meduri, Ben Trokenheim, Ricardo Urena

```
library(tidyverse)
library(tidymodels)
# add other packages as needed
```

```
wildfire <- read_csv("data/FW_Veg_Rem_Combined.csv")
```

## Introduction

Wildfires are among the most destructive natural disasters, causing significant environmental, economic, and human losses. Just this past year, California has experienced some of the most severe and costly wildfires in history, with fires destroying millions of acres, displacing thousands of residents, and costing billions in damages. The increasing frequency and intensity of these fires demonstrates the urgent need for improved wildfire response strategies. As climate change exacerbates drought conditions and fuel availability, accurately predicting fire containment time is becoming more critical for firefighting agencies and policymakers.

In this project, we will analyze wildfire characteristics, environmental conditions, and geographical factors to predict the time required to contain a fire. By using historical wildfire data, weather patterns, and vegetation information, we hope to identify key determinants of fire containment time and develop a predictive model to aid decision-making in fire management.

### Research Question:

What factors influence the time it takes to put out a wildfire, and can we develop a predictive model to estimate wildfire containment time based on environmental and fire-specific variables?

### Motivation & Importance

Wildfire response teams must quickly assess fire severity and determine the best course of action to contain and extinguish fires efficiently. However, limited firefighting resources often require prioritization based on containment difficulty. If containment time can be accurately predicted

based on fire size, vegetation type, weather conditions, and remoteness, fire departments and policymakers can make more informed decisions about resource allocation/risk mitigation.

The relevance of this research is especially important given the recent wildfires California. By understanding the key factors that impact containment time, fire agencies can allocate resources more effectively and reduce response times, which could potentially reduce the destruction caused by wildfires.

Beyond California, predicting containment time is useful for national and global wildfire management and could help with disaster preparedness and policy development efforts. By identifying the most influential factors in fire containment, we hope to provide insights to improve containment planning.

### Hypothesis

We hypothesize the following relationships between key variables and wildfire containment time:

- Larger fires will take longer to contain
- Weather conditions significantly impact containment time. Higher temperatures, stronger winds, and lower humidity levels will be associated with longer containment times, while increased precipitation will reduce containment time
- Vegetation type influences containment difficulty. Fires in densely forested areas may take longer to contain than those in grasslands or shrublands (due to fuel availability).
- Remoteness increases containment time. Fires located farther from cities and firefighting infrastructure will have longer putout times due to delayed response efforts.

### **Data description**

- The source of the data set
  - There are 4 sources of this data set:
  - [1] Short, Karen C. 2017. Spatial wildfire occurrence data for the United States, 1992-2015 [FPA\_FOD\_20170508]. 4th Edition. Fort Collins, CO: Forest Service Research Data Archive. <https://doi.org/10.2737/RDS-2013-0009.4>
  - [2] NOAA National Centers for Environmental Information (2001): Integrated Surface Hourly [1992-2015] - <ftp://ftp.ncdc.noaa.gov/pub/data/noaa/>
  - [3] Meiyappan, Prasanth, and Atul K. Jain. “Three distinct global estimates of historical land-cover change and land-use conversions for over 200 years.” *Frontiers of Earth Science* 6.2 (2012): 122-139.
  - [4] “World Cities Database.” Simplemaps, [simplemaps.com/data/world-cities](http://simplemaps.com/data/world-cities).

- A description of when and how the data were originally collected (by the original data curator, not necessarily how you found the data)
  1. Short, Karen C. (2017) - Spatial Wildfire Occurrence Data (1992-2015) When: 1992-2015  
 How: The dataset compiles wildfire occurrence records across the United States using multiple sources, including federal, state, and local fire reporting systems. It integrates data from the Fire Program Analysis (FPA) Fire-Occurrence Database (FOD), which consolidates information from agencies such as the U.S. Forest Service, Bureau of Land Management, and National Park Service.
  2. NOAA National Centers for Environmental Information (2001) - Integrated Surface Hourly (1992-2015)  
 When: 1992-2015  
 How: NOAA gathers Integrated Surface Hourly data from thousands of global weather stations, including those managed by the National Weather Service, the Federal Aviation Administration, and international meteorological organizations. The data include hourly observations of temperature, wind speed, humidity, and other meteorological variables, primarily collected through automated and manual weather stations.
  3. Meiyappan & Jain (2012) - Global Land-Cover Change & Land-Use Conversions  
 When: Datasets in the paper span from 1765 to 2010.  
 How: They used these three studies:
    1. HYDE 3.1 (Historical Database of the Global Environment): Developed by Klein Goldewijk et al. (2011), this dataset provides historical estimates of cropland and pastureland areas.
    2. RF Data: Updated estimates based on Ramankutty and Foley (1999), offering alternative historical reconstructions of agricultural land use.
    3. HH Data: Regional estimates from Houghton (2008), focusing on historical land-use changes.
  4. World Cities Database - Simplemaps  
 When: Dataset updated Regularly and there is no given timeframe.  
 How: Simplemaps compiles city location data from various authoritative sources, including government databases, geographic surveys, and other publicly available records. The dataset includes information such as city names, coordinates, population estimates, and administrative divisions.
- A description of the observations and general characteristics being measured
  - The “U.S. Wildfire Data (Plus Other Attributes)” dataset on Kaggle is a curated subset of a larger collection encompassing 1.88 million U.S. wildfire records.

Observations and General Characteristics Measured:

- Fire Occurrence Details: Each entry includes information about individual wildfire events, such as the date and location of occurrence.
- Geographic Attributes: Data points are associated with specific geographic coordinates, allowing for spatial analysis of wildfire distribution across the United States.
- Scientific Measurements: The dataset integrates supplementary scientific measurement data related to each wildfire event, such as temperature, wind speed, and humidity.

## Exploratory data analysis

<https://www.statology.org/r-extract-number-from-string/>

```
wildfire_copy <- wildfire

wildfire_copy <- wildfire_copy |>
  mutate(
    putout_time_num = parse_number(putout_time)
  )

is.numeric(wildfire_copy$putout_time_num)
```

```
[1] TRUE
```

## Analysis approach

Since the goal of our project is to understand how different fire characteristics and weather, geographic, and location-based factors influence fire containment time, `putout_time` will be the response variable in our analysis. So far, we have identified the following variables as candidate predictors that may help explain the variability in `putout_time`:

Numerical variables:

- Fire Size (`fire_size`): The total area burned by the fire (acres). We think this would be a strong predictor because larger fires generally take more time to contain (Fires with a greater burned area require more resources and extended efforts from firefighting teams.)

- **Fire Intensity (`fire_mag`):** A scaled measure of fire intensity based on fire size. We think this would be a strong predictor because more intense fires may be harder to control, and higher fire intensity might be correlated with faster fire spread and greater difficulty in containment efforts.
- **Temperature on Containment Day (`Temp_cont`):** The temperature in degrees Celsius at the fire location on the day the fire was contained. We think this would be a strong predictor because higher temperatures can dry out vegetation and increase fire spread, making containment more difficult. Lower temperatures might help slow fire progression.
- **Wind Speed on Containment Day (`Wind_cont`):** The wind speed in meters per second at the fire location on the day the fire was contained. We think this would be a strong predictor because stronger winds can spread fires rapidly and therefore make containment more challenging. Fires in areas with low wind speeds may be easier to manage.
- **Precipitation 7 Days Before Fire Containment (`Prec_pre_7`):** The amount of precipitation in millimeters at the fire location in the seven days leading up to containment. We think this would be a strong predictor because recent precipitation can increase soil and vegetation moisture, which could reduce fire intensity and help with containment efforts.

Categorical variables:

- **Cause of Fire (`stat_cause_descr`):** The reported cause of the fire (e.g., lightning, human activity, or equipment use). We think this would be a strong predictor because fires caused by natural events like lightning might have different containment challenges compared to those caused by human activities (i.e., some causes may be associated with fires that spread more aggressively).
- **Dominant Vegetation Type (`Vegetation`):** The main form of vegetation present in the fire-affected area (categorized into forest, shrubland, grassland, or urban land). We think this would be a strong predictor because different vegetation types burn at different rates; for example, denser forests/vegetation may sustain fires for longer periods of time, while grasslands might allow for quicker containment.
- **Remoteness (`remoteness`):** A non-dimensional measure of the distance to the closest city.  
We think this would be a strong predictor because fires in remote areas may take longer to contain because of limited access to firefighting resources and difficulty in transporting equipment (delayed response times).

To understand the relationship between each of our candidate predictors and fire containment time, we will first conduct EDA. For numerical predictors, we will use correlation coefficients to assess their strength of association with fire containment time. For categorical predictors, we will analyze differences in fire containment time across categories by fitting a linear regression model with categorical variables as predictors and examining the estimated coefficients,

and also rely on visualizations to determine whether there are significant differences in fire containment time across different levels of categorical variables (e.g., side-by-side boxplots to compare the distribution of fire containment time across different categories).

Based on the strongest predictors identified in our exploratory analysis, we will fit multiple linear regression models to predict fire containment time using a combination of the selected variables. Multiple linear regression is the most appropriate modeling approach for this project because our response variable (fire containment time) is a continuous, numerical variable; since we are not predicting a categorical outcome, logistic regression would not be applicable. As part of our modeling methodology, we will fit both main effects models and interaction effects models (Interactions will be included based on factors that may have a different relationship with fire containment time depending on another category based on EDA – for example, fire size may have a different effect on fire containment time depending on vegetation type).

We will fit models with 95% confidence intervals for coefficients to estimate the range in which the true slope of the relationship between each predictor and fire containment time is likely to fall. After fitting our models, we will use model assessment metrics to determine which models provide the best fit. We will evaluate models based on adjusted  $R^2$  (which accounts for the number of predictors and helps identify models that explain the most variability in fire containment time) and RMSE (average difference between observed and predicted values of the response variable). Another technique that might be useful (which we covered in class) is standardizing numerical variables in our model so we can compare each of the numerical variables and determine which are the most “impactful” predictors of fire containment time.

Based on the performance of each of our models, we will be able to draw conclusions about which factors most strongly influence fire containment time and how different environmental and fire-related conditions contribute to variability in putout time.

## Data dictionary

The data dictionary can be found [here](#) [Update the link and remove this note!]