

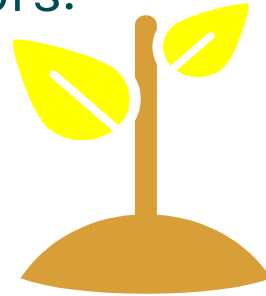
# California Wildfire Damage Projection

Egan Bailey, Sean Gu, Ani Karenovna K. Douglas Nguyen, Amanda Jo Russell

# Background

Wildfires in California are becoming increasingly common due to a variety of factors:

- Climate change
- Drought Frequency & Seasonal Rains
- Human Causal Factors
- Historical Forest Mismanagement
- Increased population in vulnerable areas



## Background (Cont.)

FEMA would benefit from a model projects the costs resulting from a fire disaster.

- ◆ Financial aid can then be properly allocated to the impacted communities.



# Data Science Problem

**Our Goal is to Predict the Dollar Amount of  
Damage done to specific counties within the state of  
California due to wildfires**



# Data Science Problem

Interaction between some of these variables have proven useful for making predictions:

1. Average size of the fire spread
2. Number of fire origins within a given area
3. Weather conditions
4. Variables representing relative county affluence
5. Percentages of wilderness/developed land in a given county

# Data Science Process / Methodology

- ◆ Data Collection
- ◆ Cleaning data & EDA
- ◆ Feature Engineering
- ◆ Model Fitting & Tuning
- ◆ Evaluations



1.

# Data Collection

# Data Collection: Sources



## ◆ Cal Fire: Redbooks (2010 - 2016)

- ◆ Number of Fires (sizes and causes)
- ◆ Dollar Damage



## ◆ Yelp

- ◆ Number of Campsites/RV Parks per County

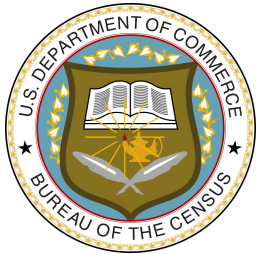


# Data Collection: Sources



## ◆ National Oceanic and Atmospheric Administration (NOAA)

- ◆ Weather Data



## ◆ National Census Bureau/Wikipedia

- ◆ Population Data

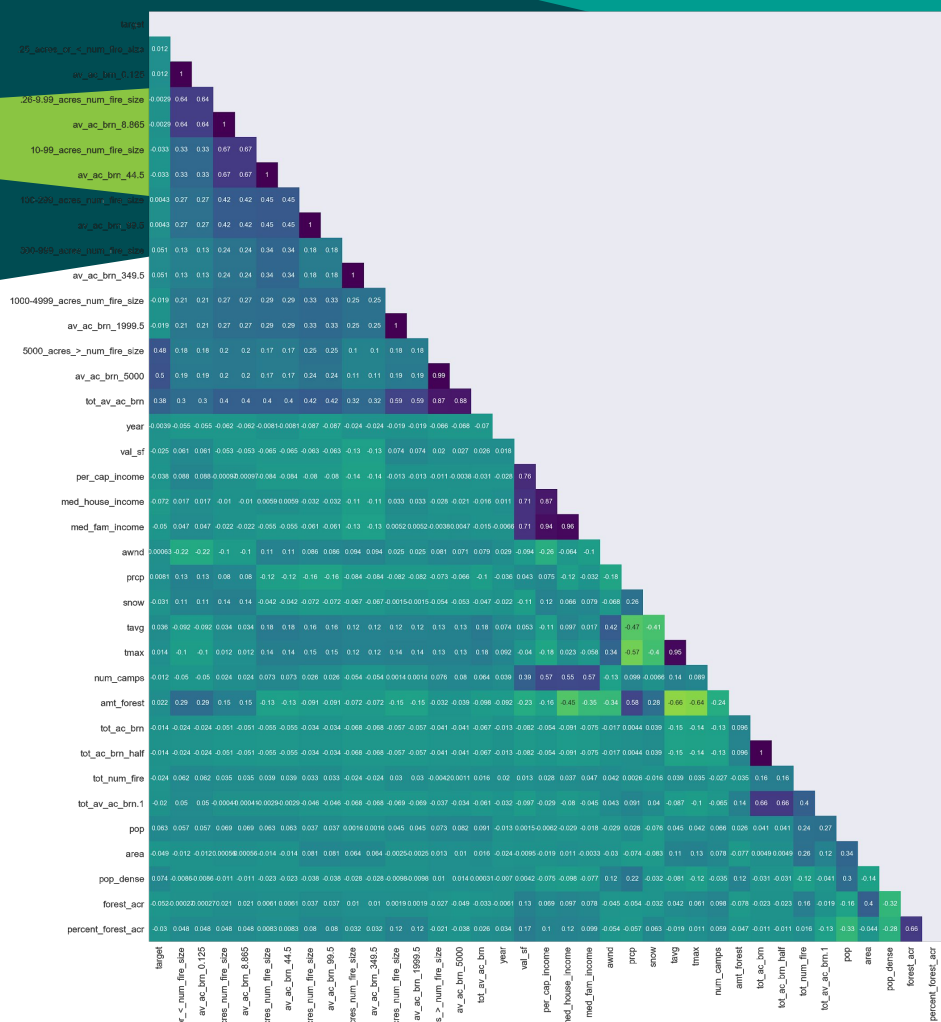
## ◆ Employment Development Department

- ◆ Workforce Data



## 2. Data Cleaning & EDA

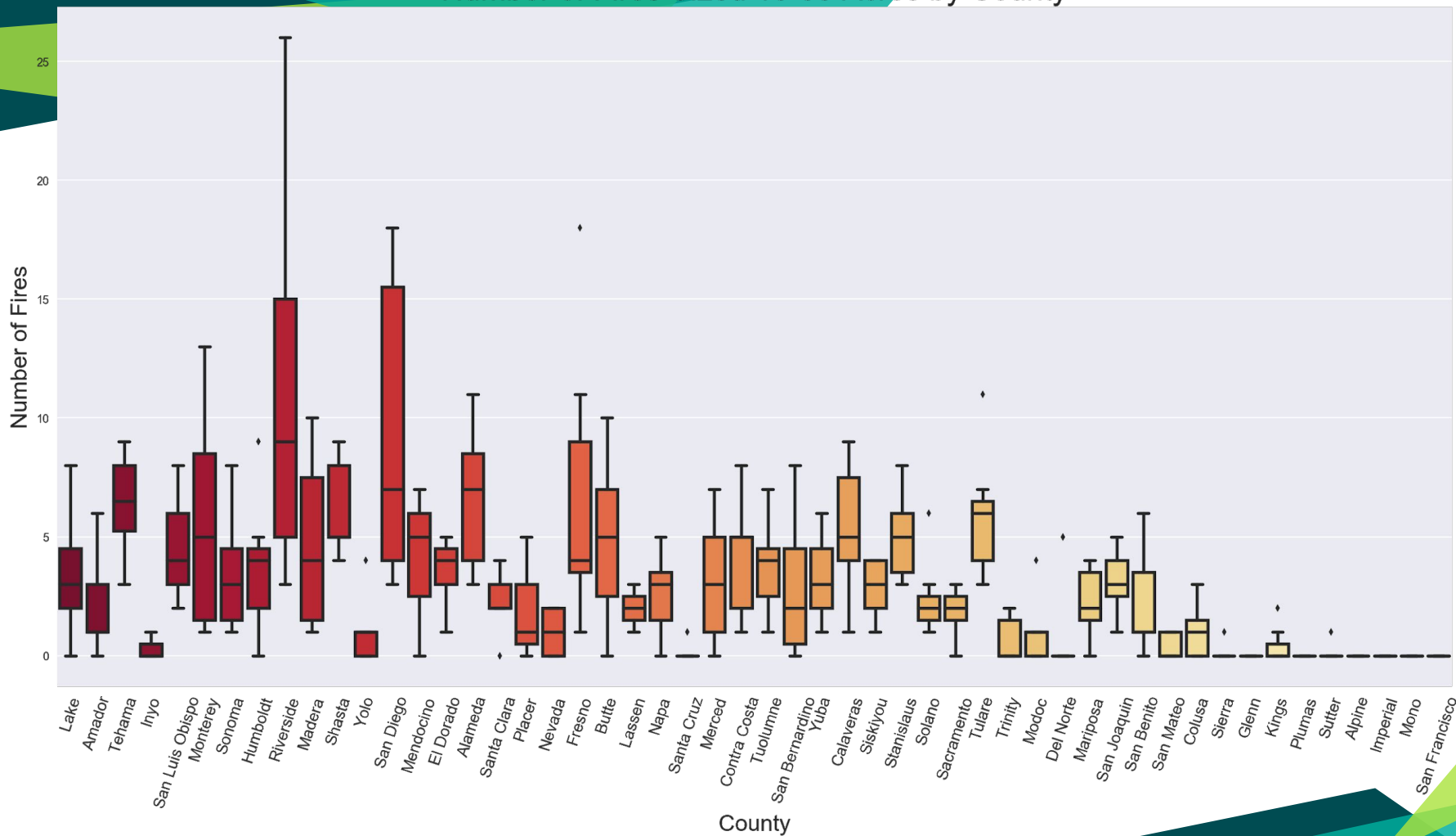
# Correlation: Heatmaps



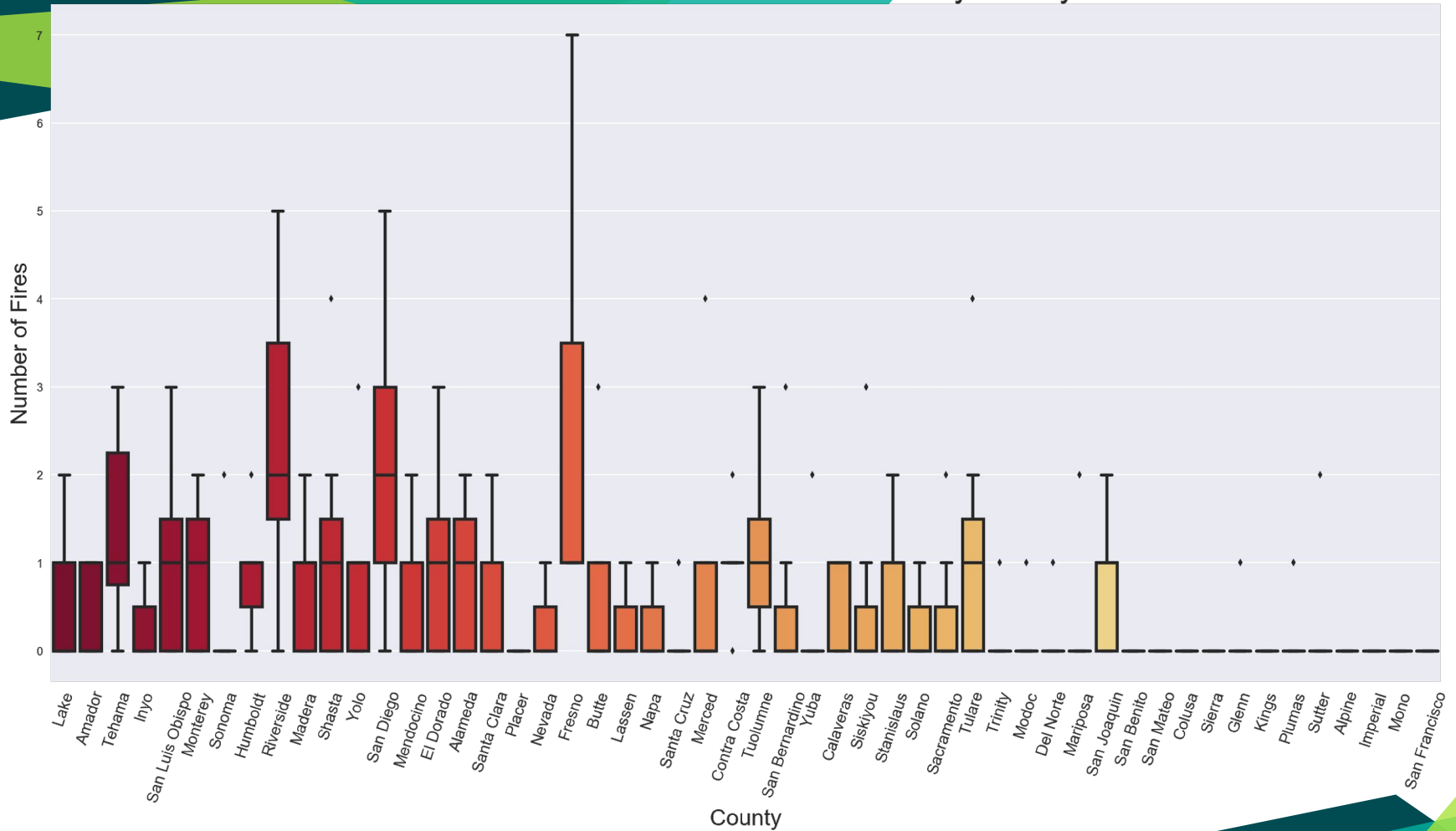
target	1
av_ac_brn_5000	0.5
5000_acres_num_fire_size	0.48
tot_av_ac_brn	0.38
pop_dense	0.074
pop	0.063
av_ac_brn_349.5	0.051
300-999_acres_num_fire_size	0.051
tavg	0.036
amt_forest	0.022
tmax	0.014
av_ac_brn_0.125	0.012
.25_acres_or_less_num_fire_size	0.012
prcp	0.0081
100-299_acres_num_fire_size	0.0043
av_ac_brn_99.5	0.0043
awnd	0.00063
av_ac_brn_8.865	-0.0029
26-9.99_acres_num_fire_size	-0.0029
year	-0.0039
num_camps	-0.012
tot_ac_brn	-0.014
tot_ac_brn_half	-0.014
av_ac_brn_1999.5	-0.019
1000-4999_acres_num_fire_size	-0.019
tot_av_ac_brn.1	-0.02
tot_num_fire	-0.024
val_sf	-0.025
percent_forest_acr	-0.03
snow	-0.031
av_ac_brn_44.5	-0.033
10-99_acres_num_fire_size	-0.033
per_cap_income	-0.038
area	-0.049
med_fam_income	-0.05
forest_acr	-0.052
med_house_income	-0.072

target

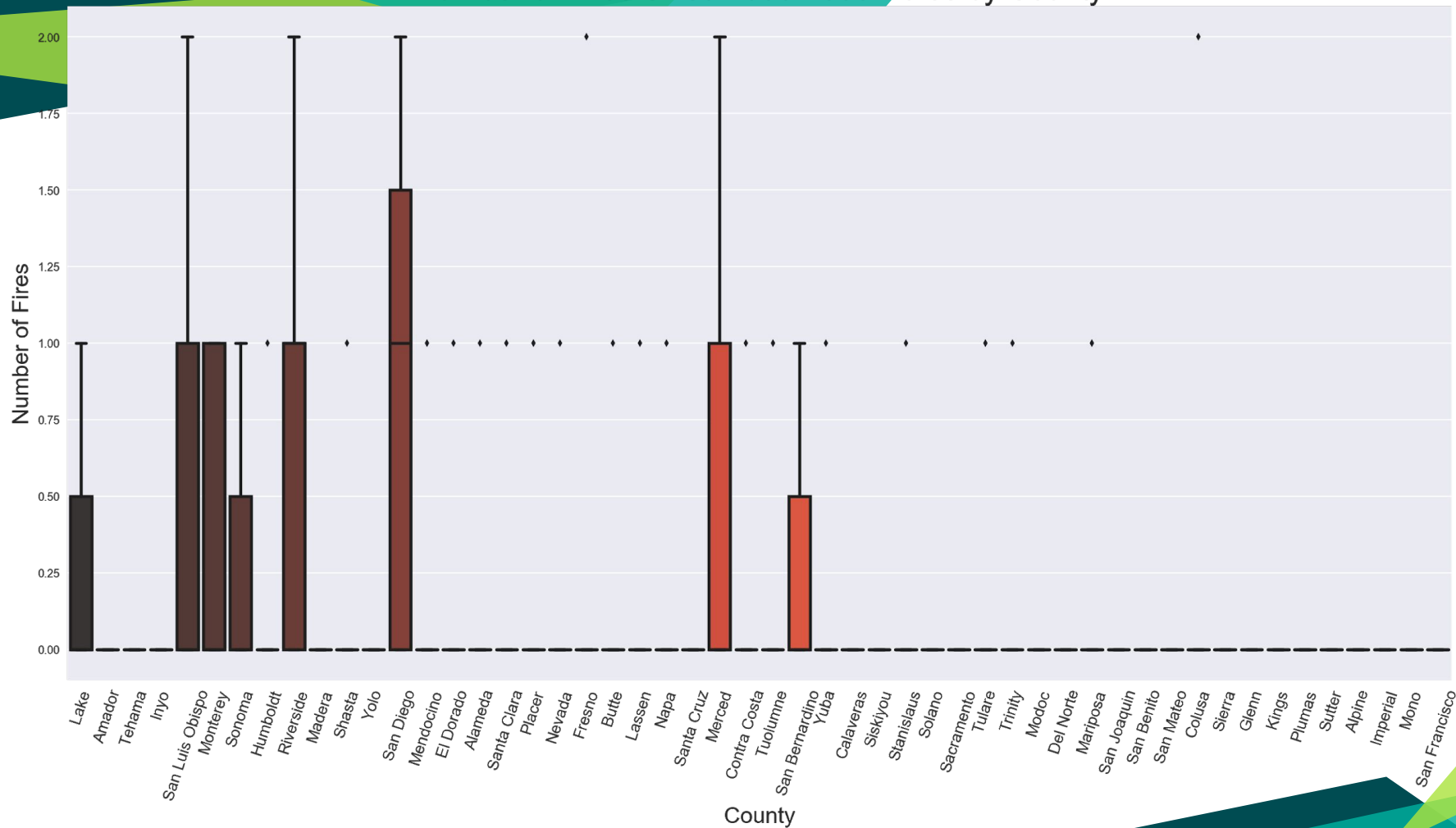
# Number of Fires Sized 10-99 Acres by County



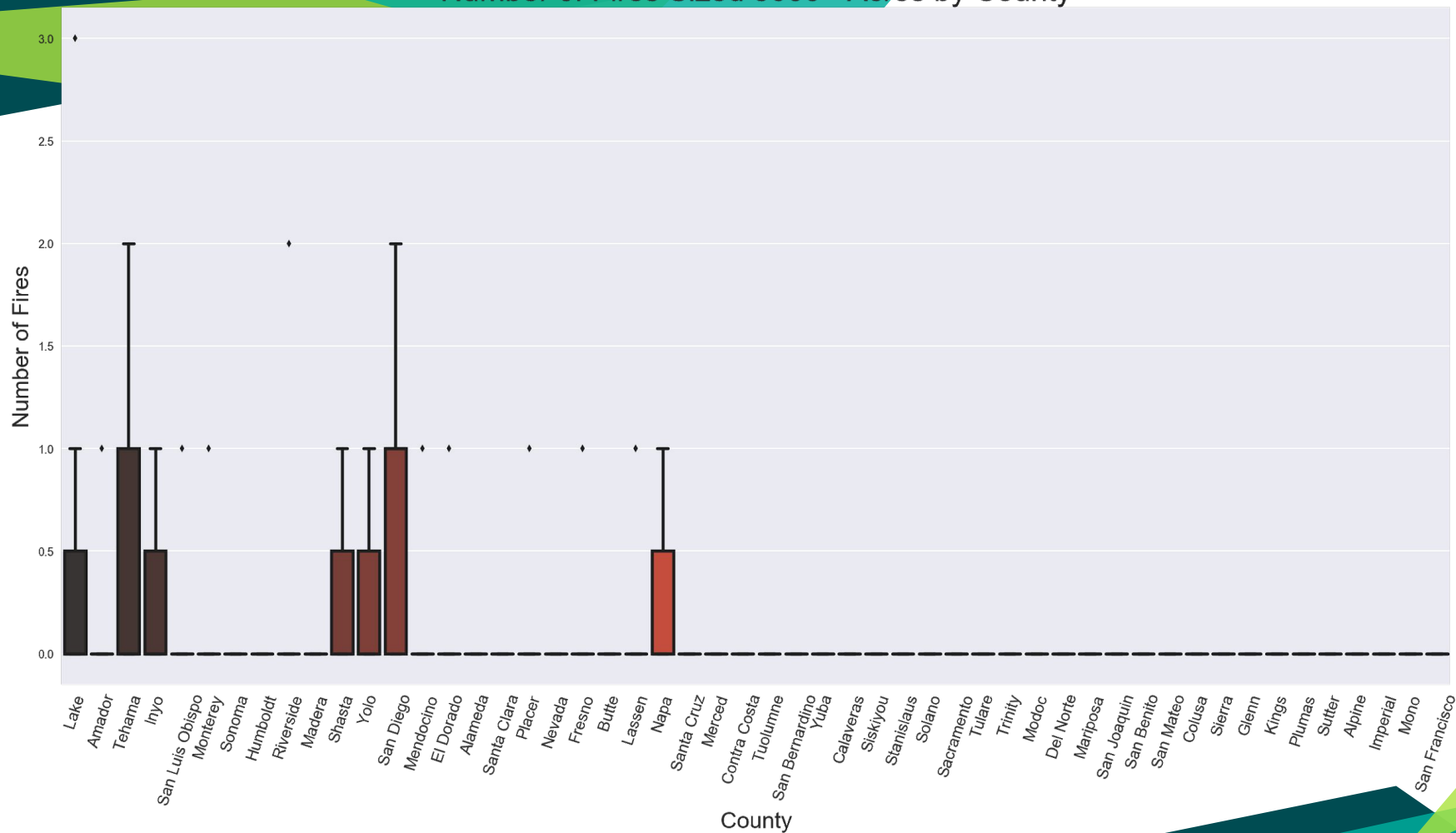
# Number of Fires Sized 100-299 Acres by County



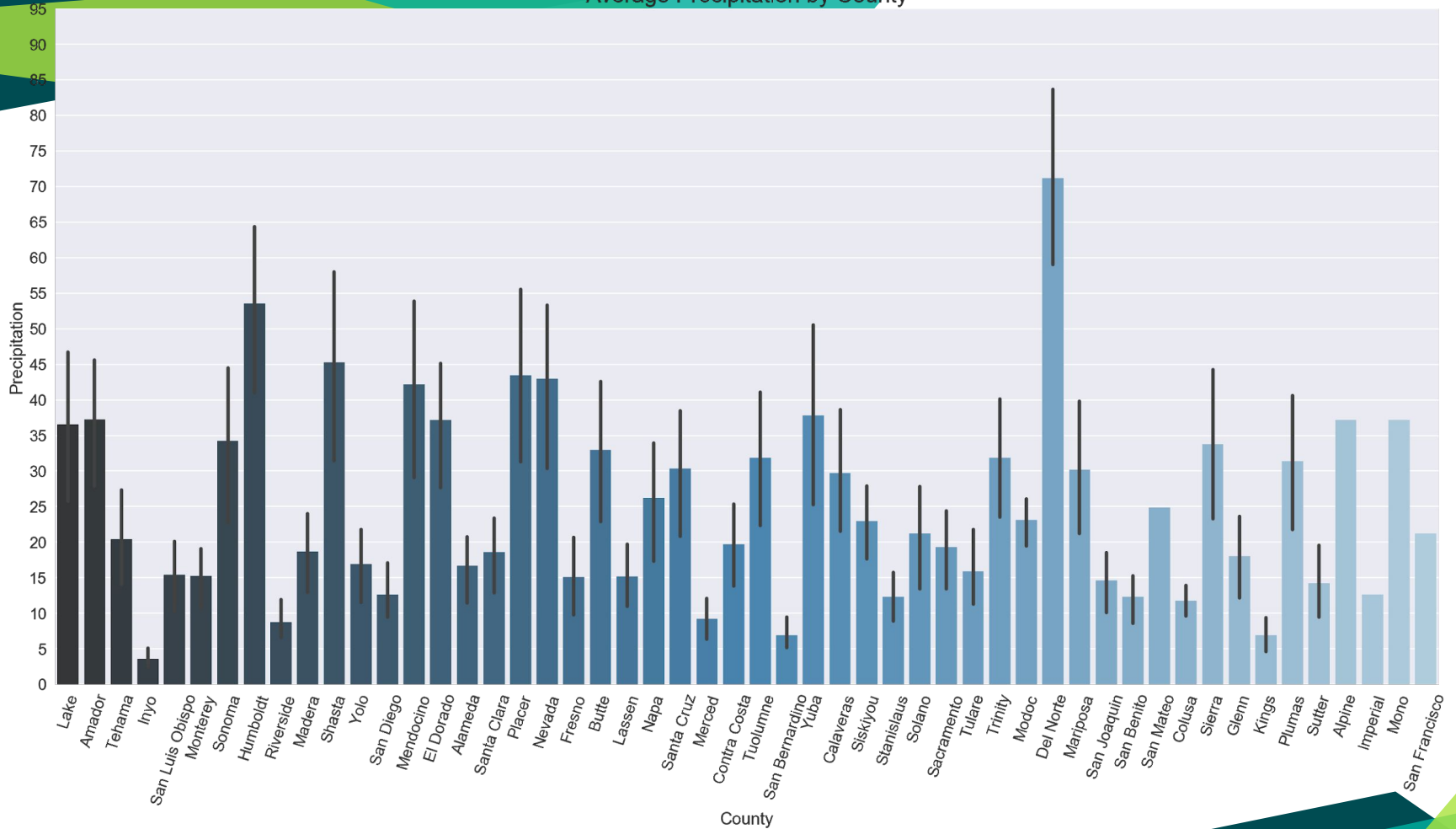
# Number of Fires Sized 1000-4999 Acres by County



## Number of Fires Sized 5000+ Acres by County

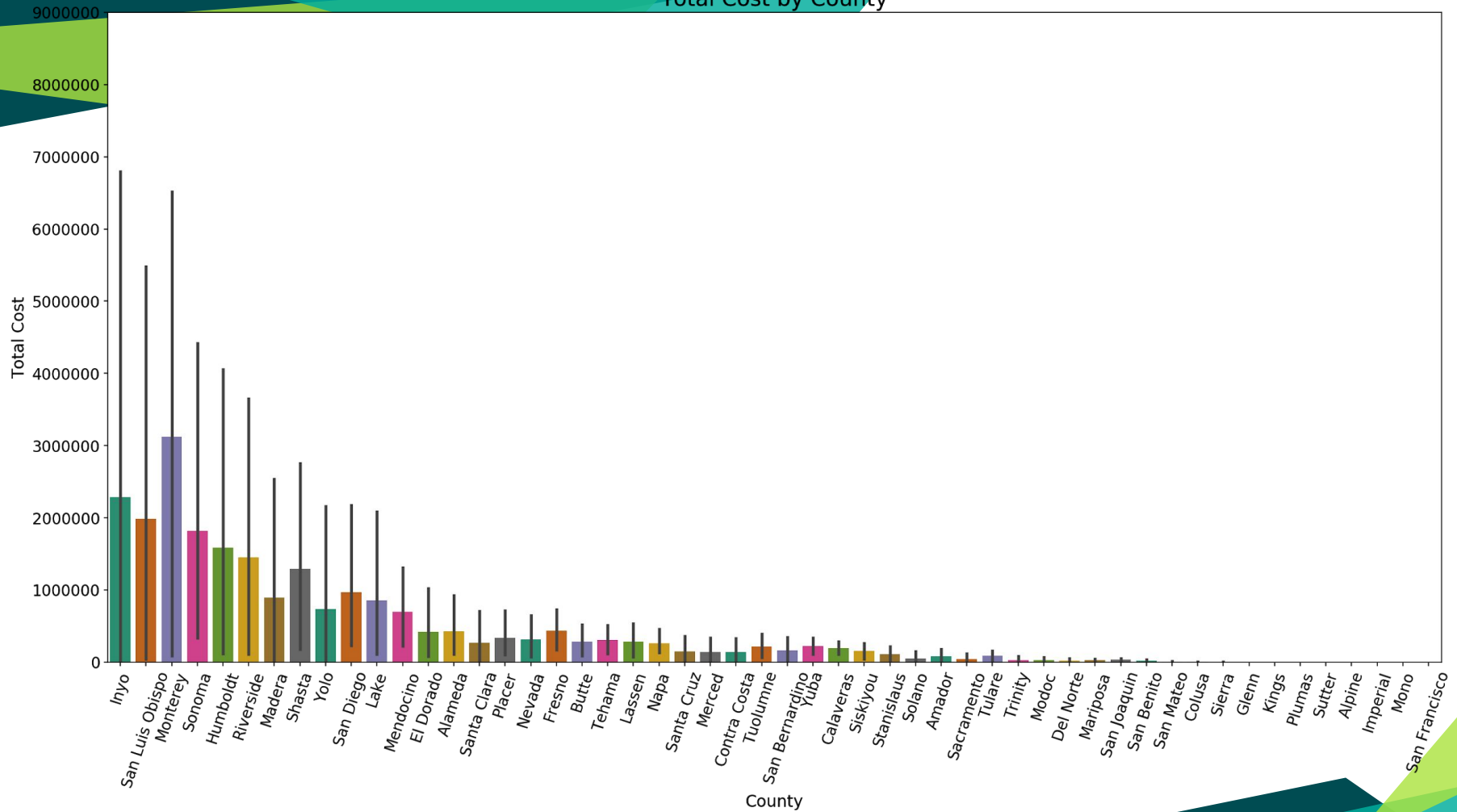


## Average Precipitation by County





## Total Cost by County





# 3. Feature Engineering



4.

# Model Fitting and Tuning

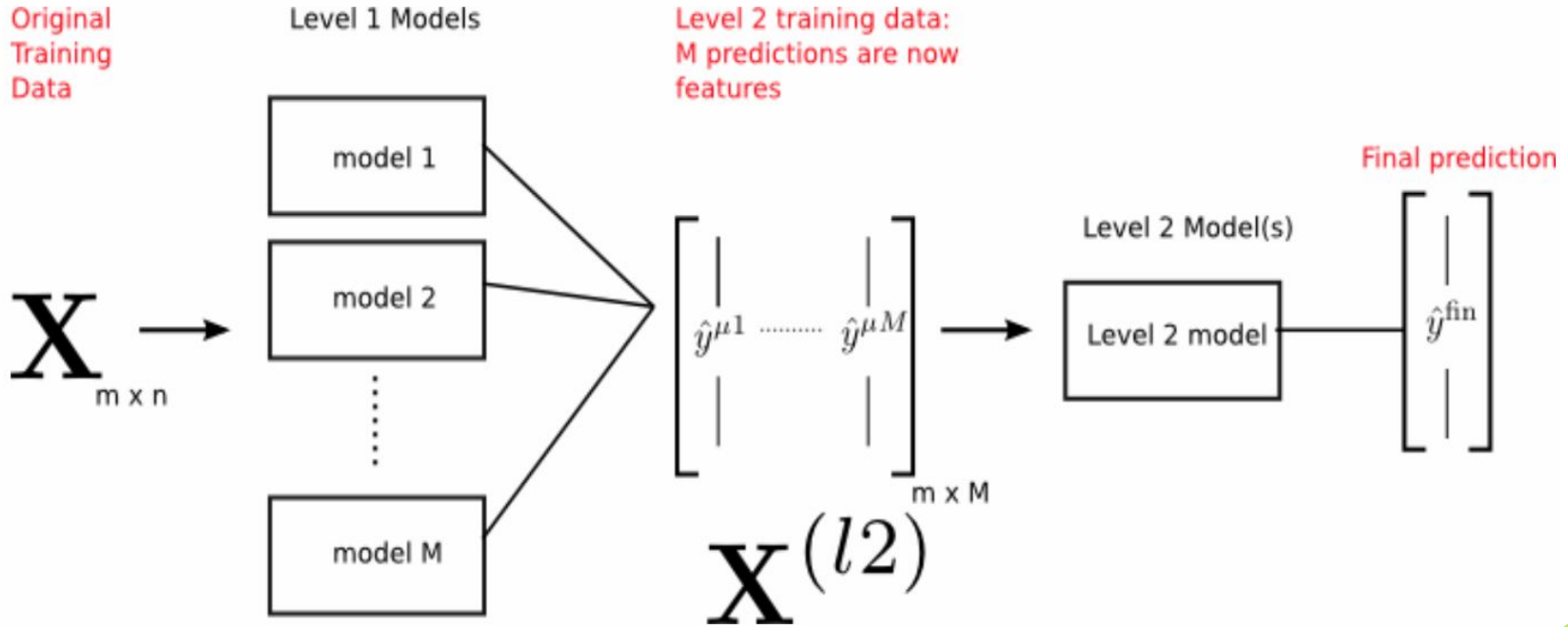
# The Models: Regression

1. Multiple Linear Regression
2. Lasso
3. Ridge
4. KNN
5. SVR
6. Random Forest
7. Bagging
8. Extra Trees
9. Ada Boost
10. Gradient Boost

Each model was trained and fit in three different contexts:

1. Base Case - no tuning
2. Grid Search - tuning
3. Stacking - predictions from best gs models become the new features

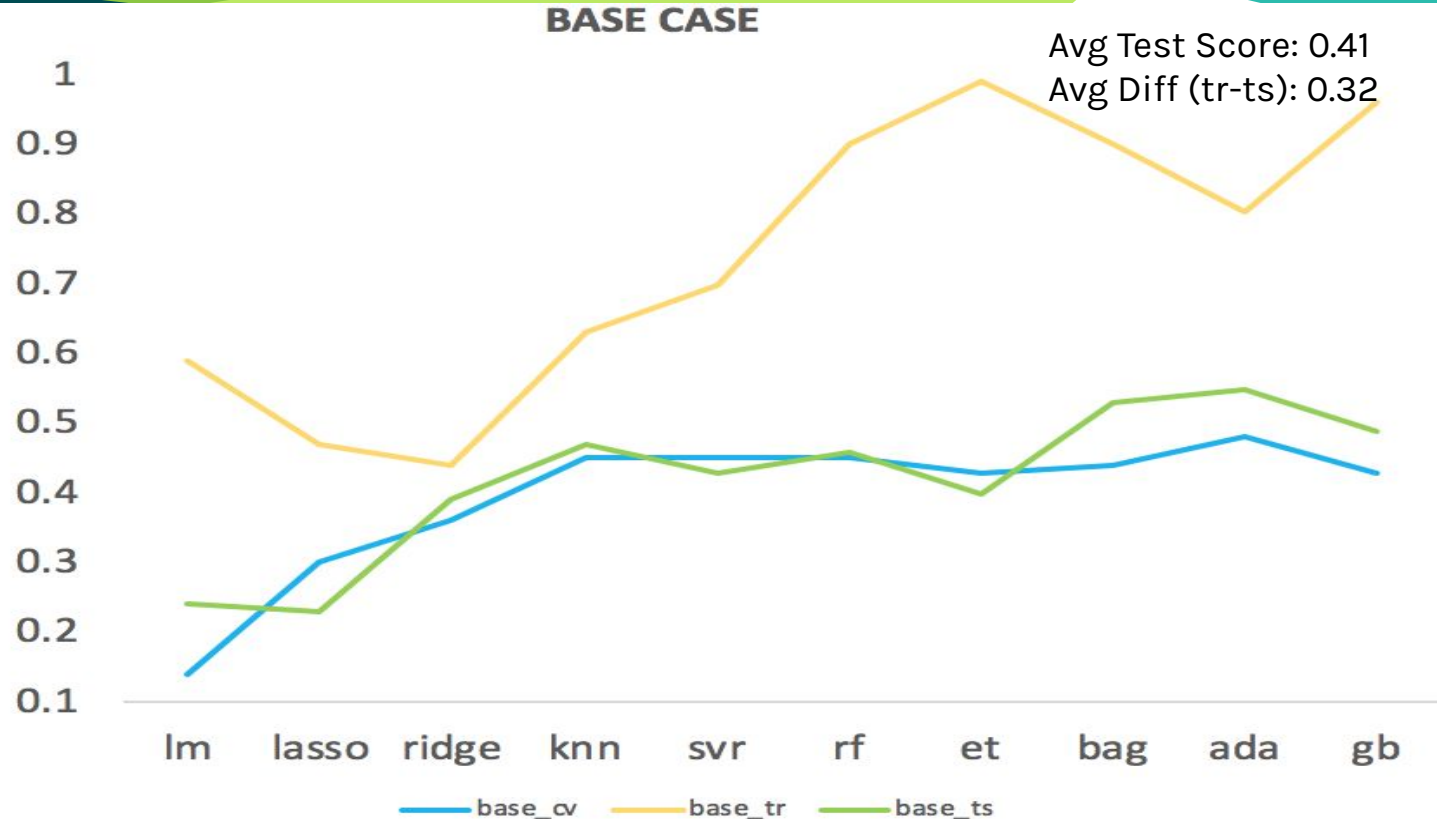
# Stacking



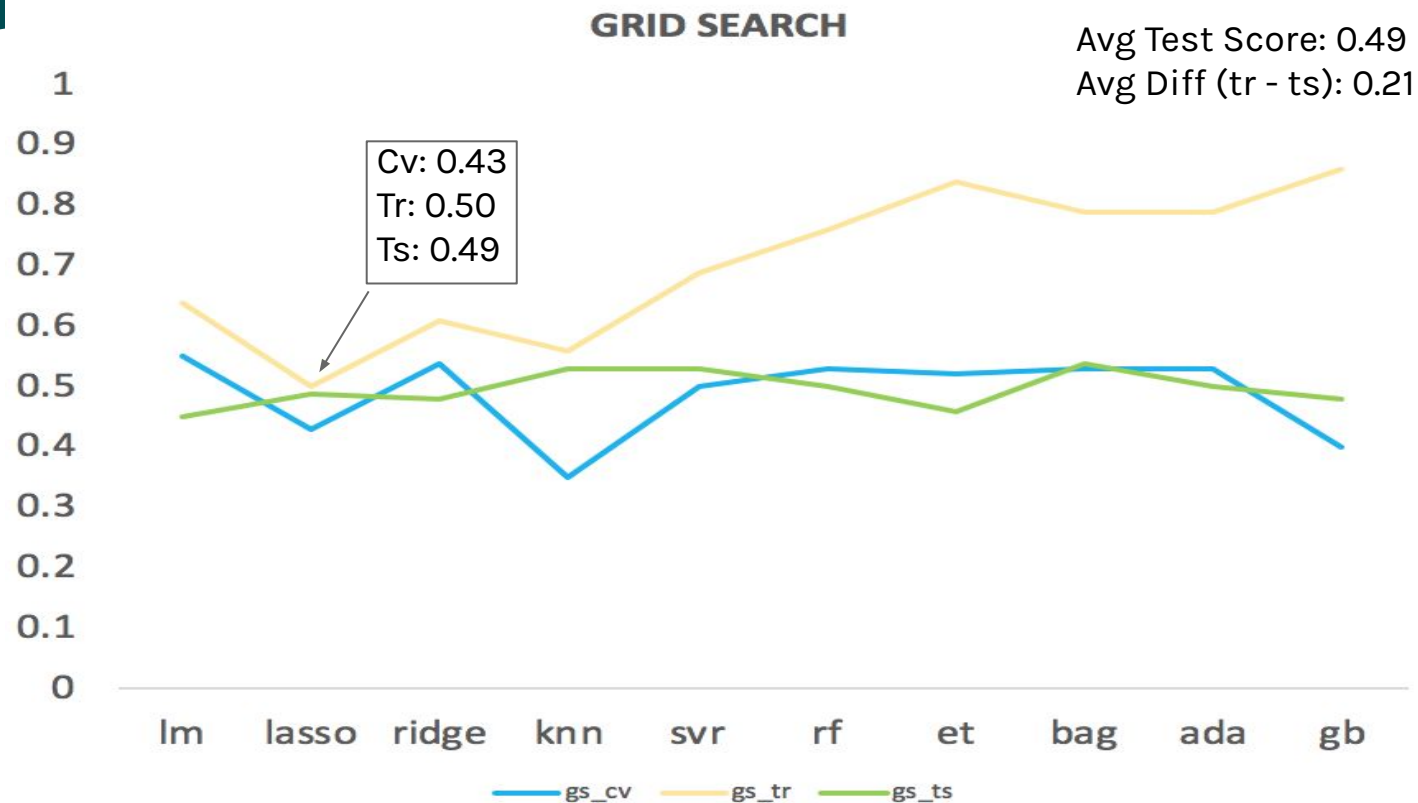


# 5. Evaluation and Conclusions

## Train Scores, Test Scores and CV Scores (R\_squared) for Each Model

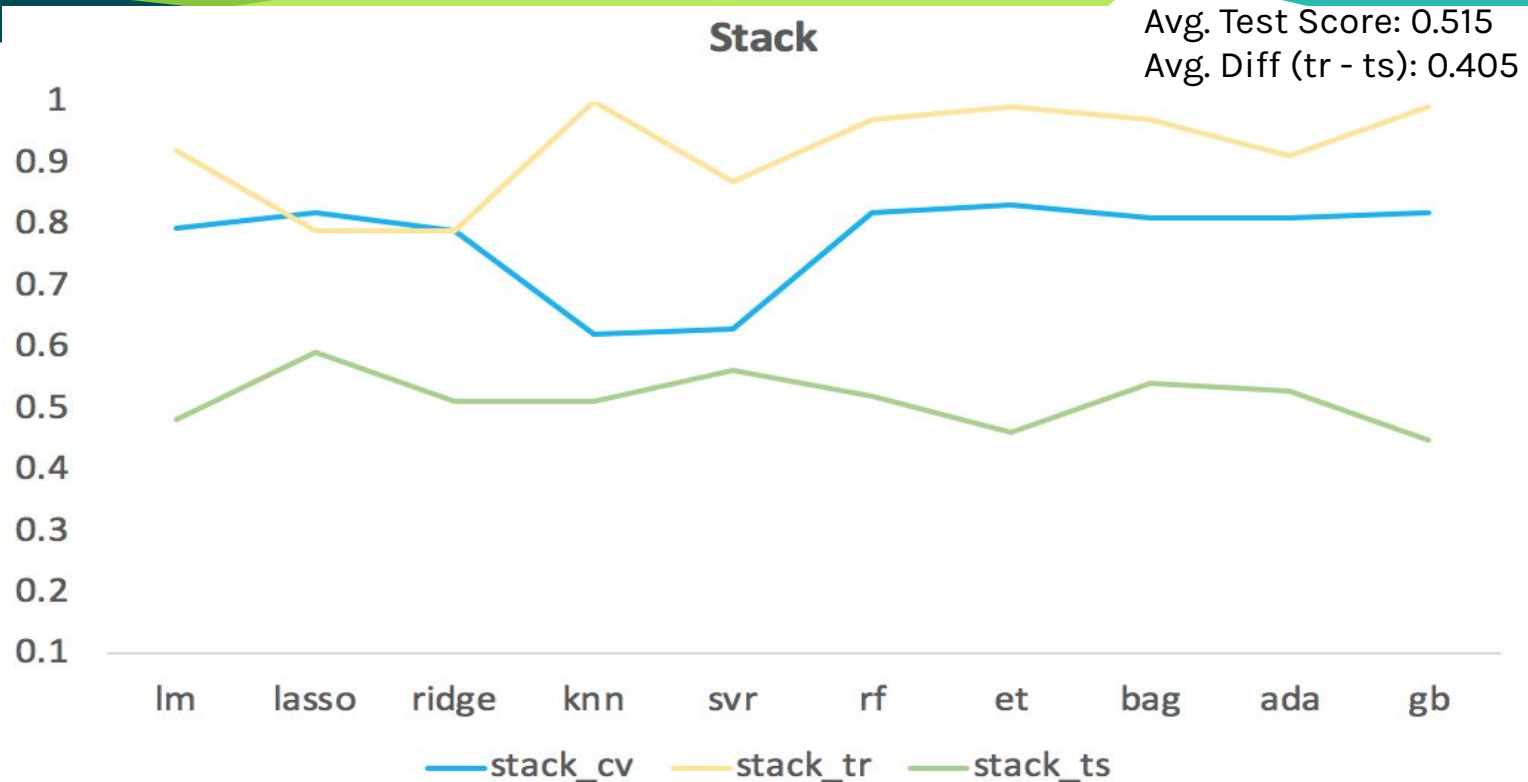


## Train Scores, Test Scores and CV Scores (R\_squared) for Each Model





## Train Scores, Test Scores and CV Scores (R\_squared)for Each Model



# Demonstration of Our Solution

Interaction between some of these variables have proven useful for making predictions:

1. Average size of the fire spread
2. Number of fire origins within a given area
3. Weather conditions
4. Variables representing relative county affluence
5. Percentages of wilderness/developed land in a given county

# How to Use the Model

- The best model is optimized for prediction, not for explainability
- Input: features
- Output: dollar amount damage

# Limitations of the Model

- Data Collection Issues
  - Accuracy and completeness of raw data
- Missing values
  - Assumptions made while imputing missing values
- Overall complexity of variable space that are predictive of fire processes

## Next Steps

- CalFire data that is more granular, so as to include detailed data on individual wildfires (e.g dollar amount damage, acres burned, etc).
- Data on the epicenter of fires to help determine spatial correlations.
- An examination of the impact of US Forest Service “controlled burns” and how it affects future fires.
- Insurance premium increases.
- Fire suppression/containment data.

# Sources

- <https://www.nytimes.com/2018/11/09/climate/why-california-fires.html>
- <https://www.kqed.org/science/1927354/controlled-burns-can-help-solve-californias-fire-problem-so-why-arent-there-more-of-them>
- [https://en.wikipedia.org/wiki/List\\_of\\_counties\\_in\\_California](https://en.wikipedia.org/wiki/List_of_counties_in_California)
- [https://www.fire.ca.gov/fire\\_protection/fire\\_protection\\_fire\\_info\\_redbooks](https://www.fire.ca.gov/fire_protection/fire_protection_fire_info_redbooks)
- <https://www.pewtrusts.org/en/research-and-analysis/reports/2018/06/19/what-we-dont-know-about-state-s-pending-on-natural-disasters-could-cost-us>
- <https://www.kdnuggets.com/2017/02/stacking-models-improved-predictions.html>
- <https://www.codecademy.com/articles/seaborn-design-ii>
- <https://www.ncdc.noaa.gov/cdo-web/search>



# Thanks!

## Any Questions?