



01 INTRODUCTION

Wildfires pose as a formidable menace to both the environment and human well-being. Despite the presence of Wildfire prediction and control centers such as European Forest Fire Information System, NIFC, and others, the global landscape continues to witness the occurrence of extensive amount of wildfires. These catastrophic events wreak havoc on ecosystems, causing irreparable damage, and tragically, leading to the loss of numerous lives.

Wildfires also tend to generate substantial volumes of finer particulates recognized as PM_{2.5}, along with even smaller nanoparticles, both of which are well known for their detrimental effects on human health.

The smoke emerging from forest and peat combustion can persist in the atmosphere for extended periods, spanning weeks, and traverse great distances of thousands of kilometers. This phenomenon adversely affects the well-being of populations residing in regions far removed from the actual wildfire sites.

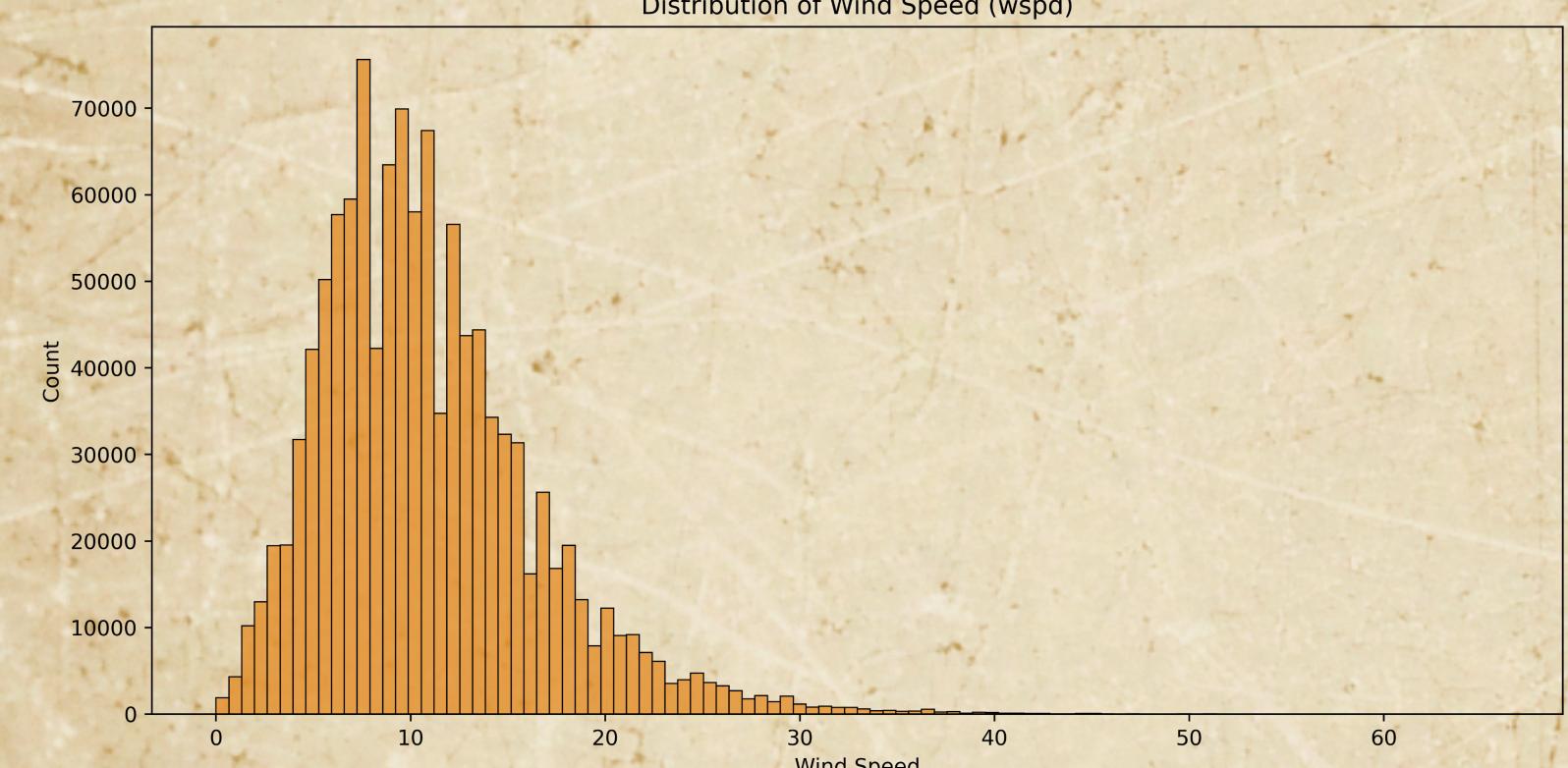
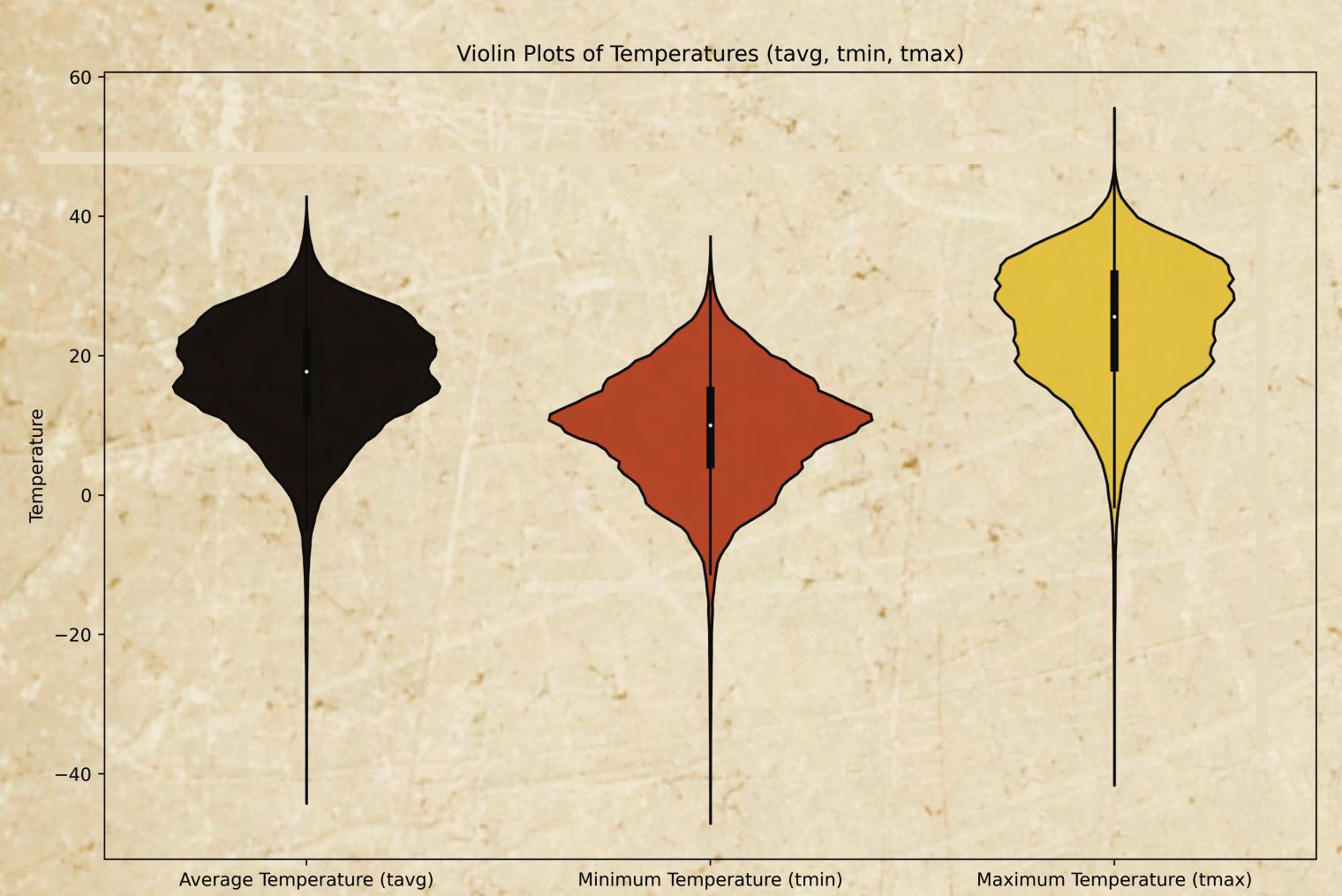
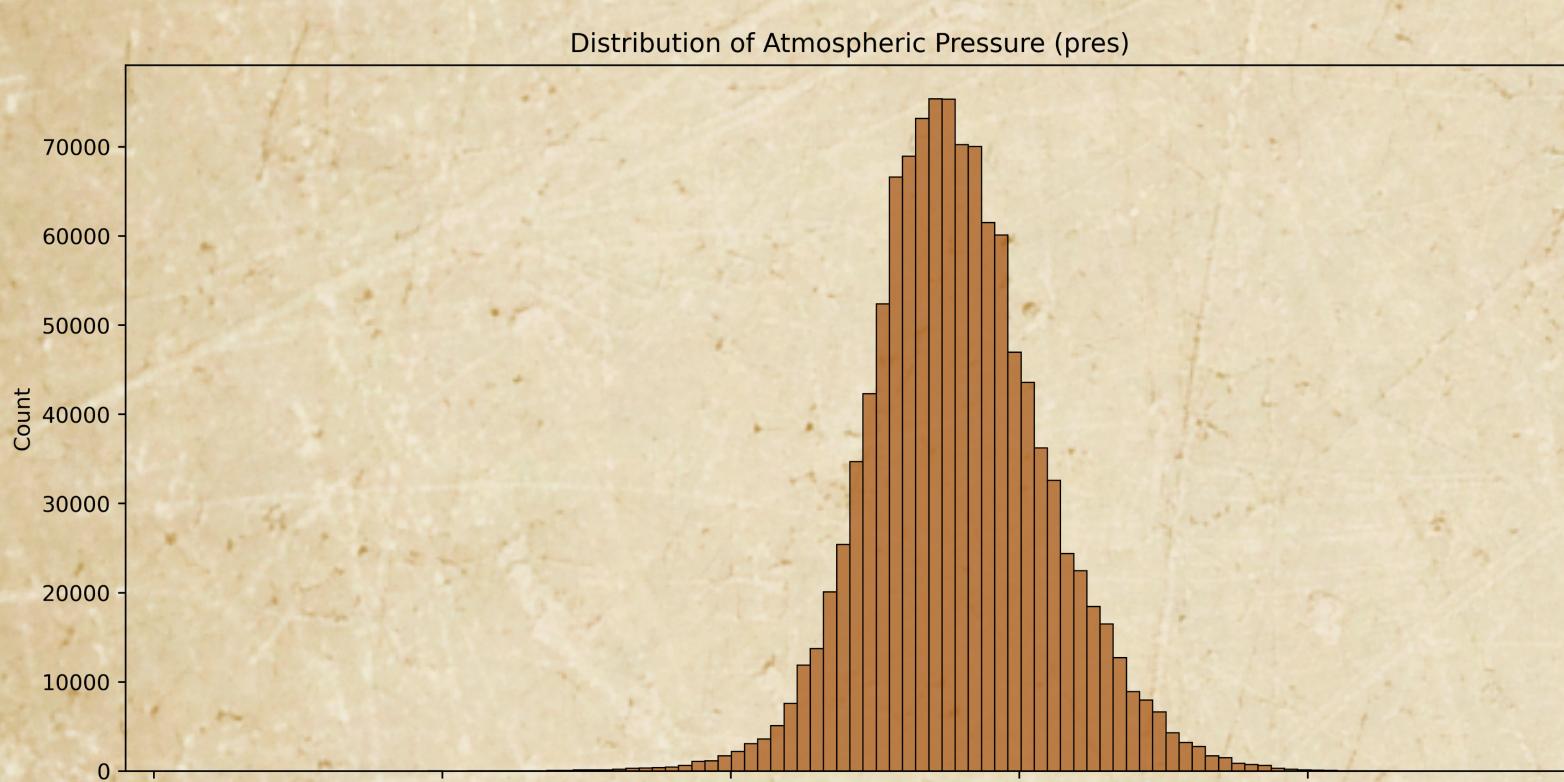
04 EXPLORATORY DATA ANALYSIS

DESCRIPTIVE STATISTICS: Upon compiling our dataset, we began with a comprehensive statistical summary. The numerical columns such as Fire Size, Latitude, Longitude, tavg, and weather parameters revealed patterns and distributions crucial for our understanding. For instance, the average fire size stood at approximately 762 hectares, while the mean atmospheric pressure was around 1015 hPa.

On the categorical front, our dataset consisted of diverse entries. The Fire Class Size column, our primary point of interest, contained 8 unique categories, and the State column highlighted wildfires spread across 51 U.S. states.

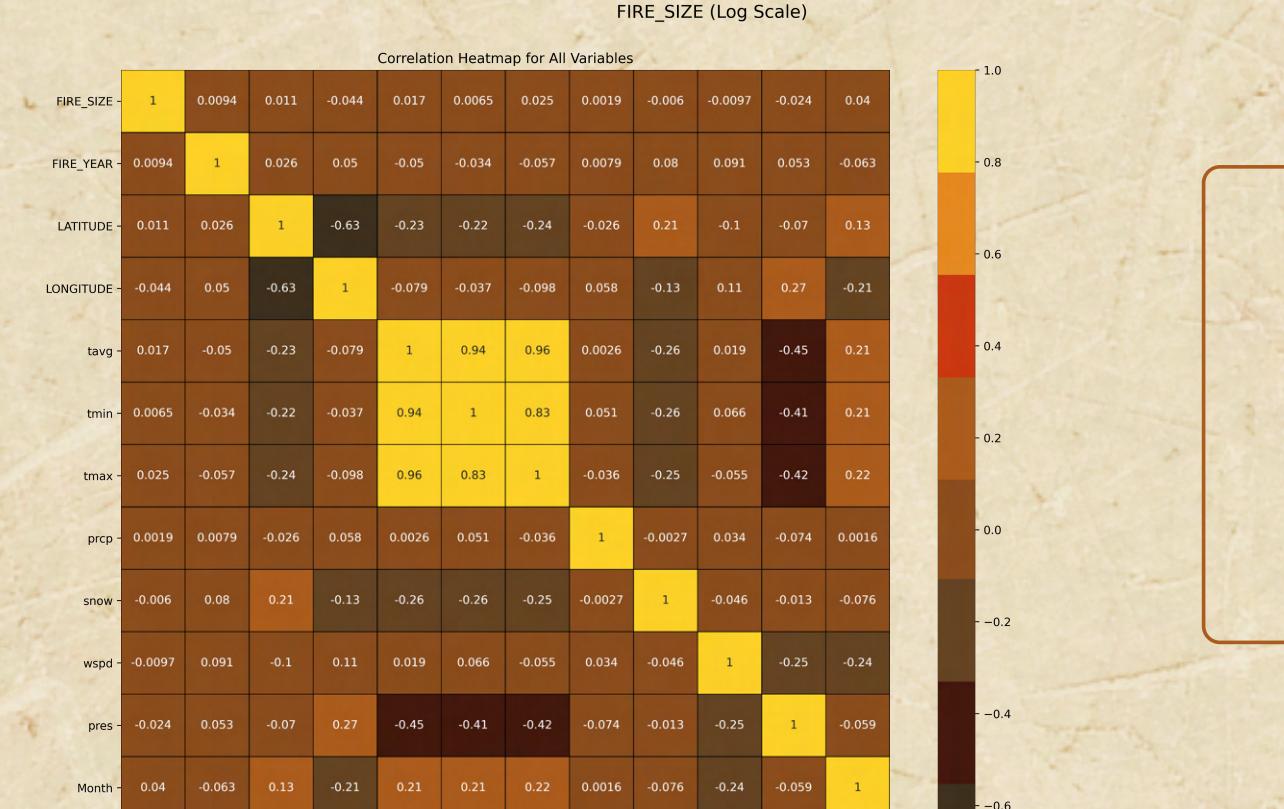
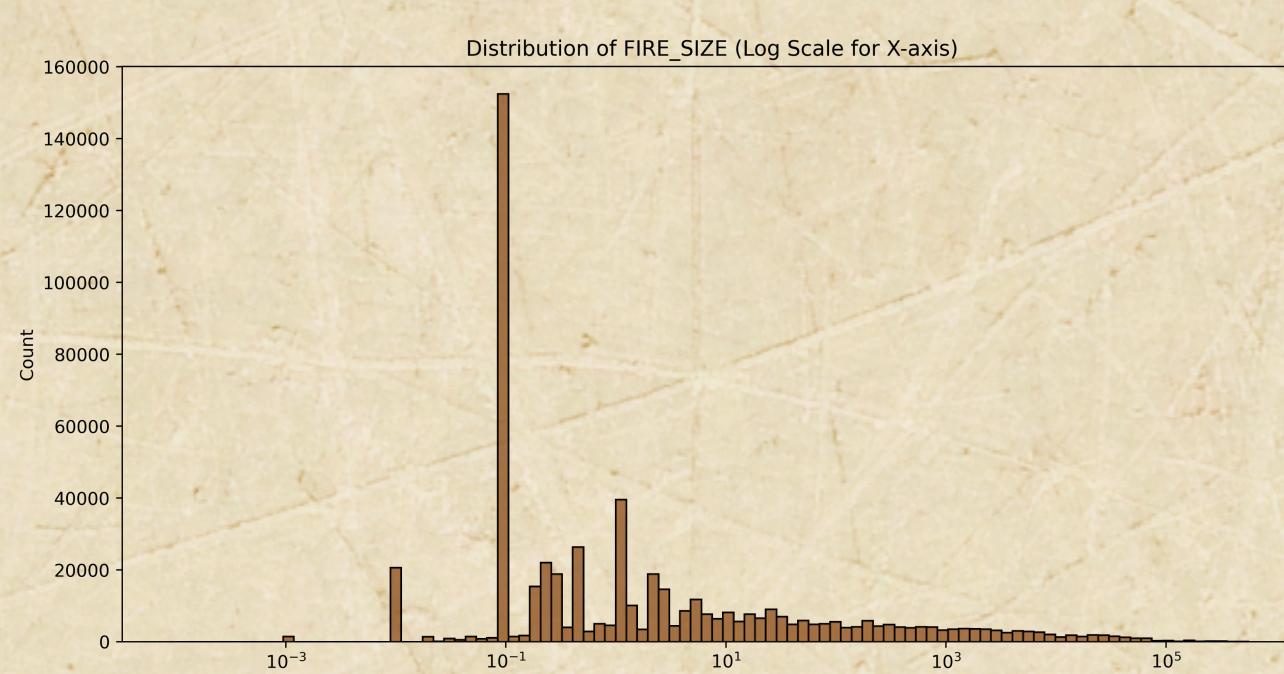
GRAPHICAL ANALYSIS: To discern patterns and relationships visually:

Weather Parameters: Histograms depicting the distributions of atmospheric pressure, wind speed, and temperature. Violin plot depicting distribution of temperature variables. This gave insights into the typical weather conditions on days with wildfires.



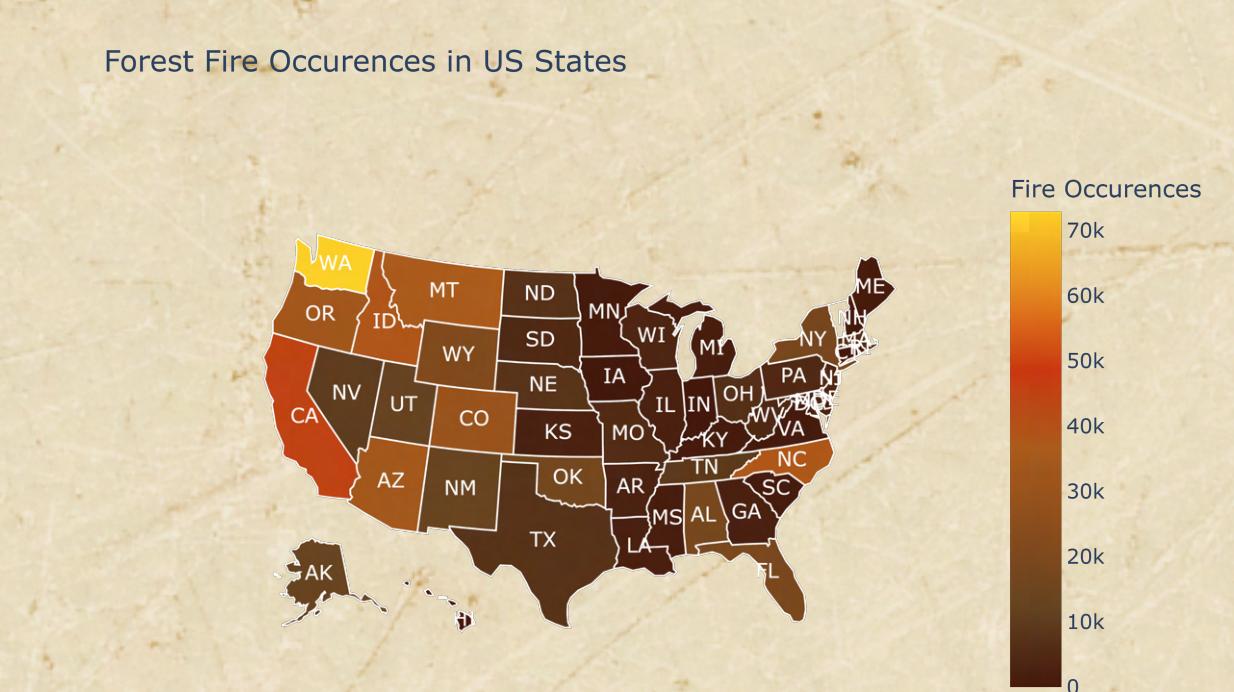
Fire Size: We took the logarithmic scale of the fire sizes to handle its wide range and plotted its distribution.

Fire Size Distribution: Lastly, the distribution of 'FIRE_SIZE_CLASS' was depicted through a pie chart, showing the relative frequencies of different fire sizes.



Correlation: A heatmap showcased the relationships between variables, illuminating potential multicollinearity and influential factors.

State-wise Comparisons: We visualised wildfire occurrences across U.S. States using a heat map, highlighting region with frequent outbreaks and areas experiencing intense fires.



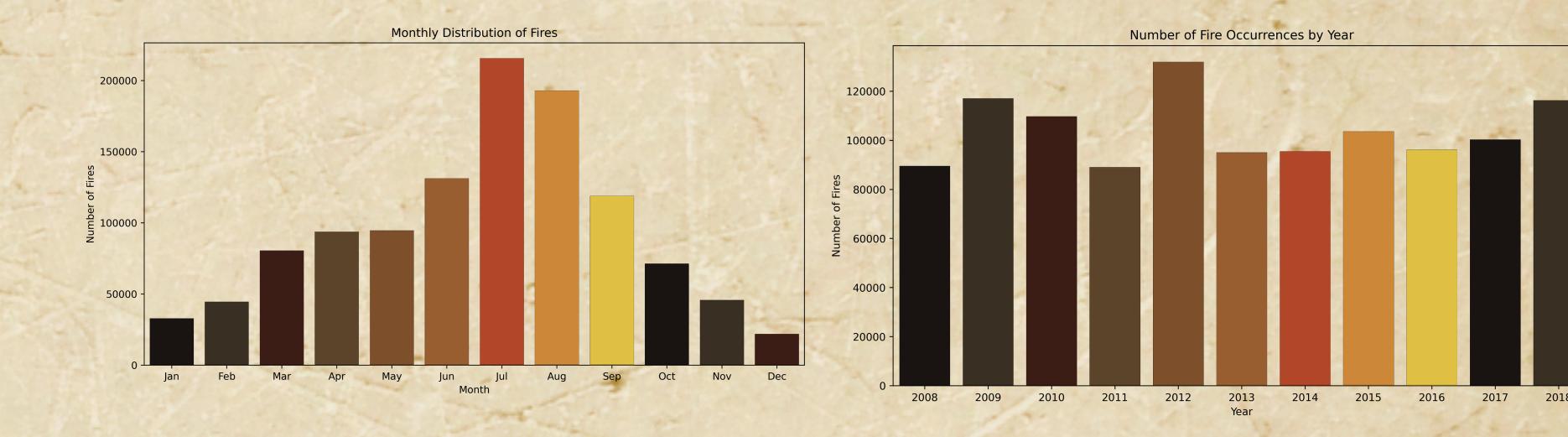
FIRE SIZE CLASS DESCRIPTION

A	B	C	D	E	F	G	N
0-25 Acres	26-9.9 Acres	10-99.9 Acres	100-299 Acres	300-999 Acres	1000-4999 Acres	5000+ Acres	No Fire

TEMPORAL ANALYSIS

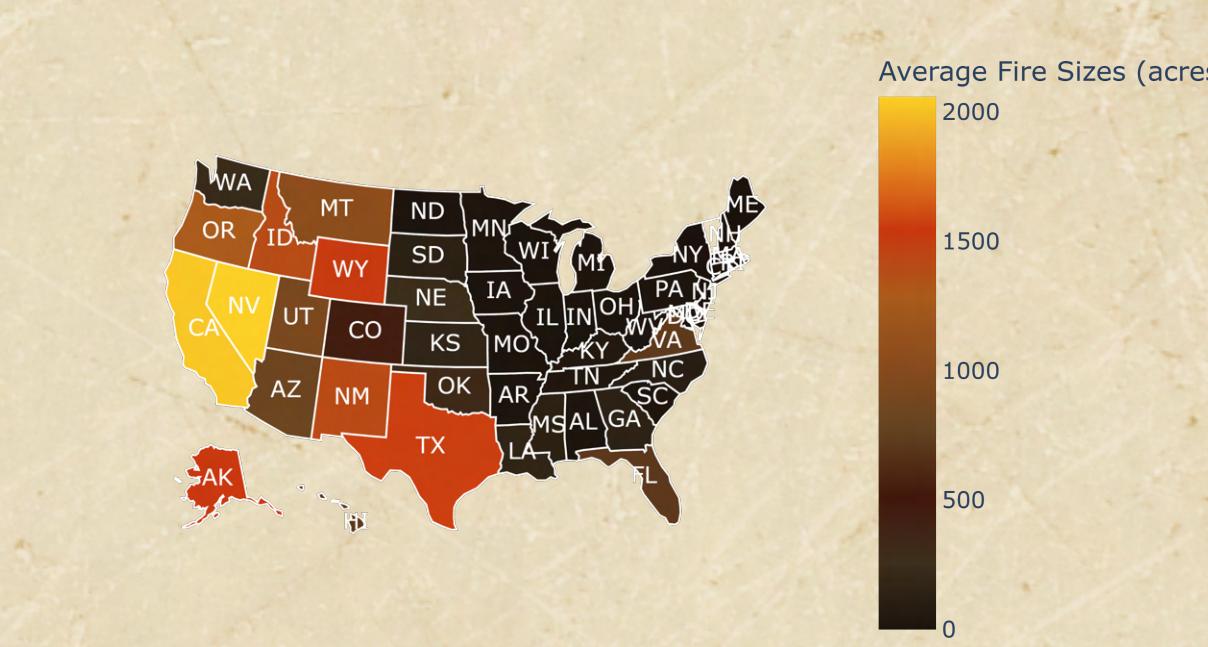
Fire Dynamics: Monthly fire distributions and their causes were charted, revealing seasonality and common ignition sources

Temporal Analysis: The yearly distribution highlighted the number of fire occurrences, emphasizing potential increasing or decreasing trends.



Cause Analysis: Employed a bar plot to describe the primary cause of wildfire revealing human activities and natural events as predominant ignition sources.

US States by average forest fire sizes



Missing Values: Our initial date range (2013-2018) left us with a diminished dataset of just over 400,000 rows after handling missing data. This constraint led us to expand to 2008-2018, yielding 1.4 million usable rows.

Computational Constraints: The sheer size of the dataset, combined with the desire to utilize more complex models like Random Forests, led to computational challenges. The balance between model performance and efficiency became a pivotal aspect. Our computational limitations could have hindered Random Forest from realizing its full potential. Issues such as byte allocation and extended training times were frequent hurdles.

Data Imbalance: Our strategy to distribute 'N' values helped address imbalance. However, omitting 'N' values showed a notable accuracy boost, underscoring the influence of balance on performance.

In essence, we grappled with challenges tied to data volume, missing values, computational limitations, and data balance, each demanding tailored solutions for effective modeling.

09 FUTURE WORK

Dataset Expansion: Our accuracy saw a direct boost with a larger dataset. This advocates for further data augmentation, perhaps extending the date range or integrating additional data sources.

Satellite Data Integration: Satellite imagery offers insights into vegetation health, soil moisture, and other crucial variables. Incorporating this can enhance the model's predictive capability.

Real-time Monitoring Systems: Implementing real-time data capture mechanisms, like sensor networks, can provide timely updates, enhancing predictive accuracy.

Climate Change Impact Analysis: A deeper dive into climate change effects and global warming trends can offer valuable insights. This knowledge can refine hyperparameter tuning, optimizing model performance.

Feedback Loop: Establishing a feedback mechanism is crucial. After predictions are made and interventions implemented, evaluating the outcomes and retraining the model with this feedback ensures its continuous evolution and improvement.

10 CONCLUSION

In our journey to predict wildfire occurrences and sizes using a decade-long dataset, we discovered the critical interplay between meteorological, geographical, and temporal factors. The unexpected success of the decision tree model and the marked improvement upon incorporating geographical data emphasize wildfires' multifaceted nature. Our work suggests a need for wildfire management to move beyond just meteorological data, incorporating a more holistic approach. The challenges, from data volume to computational limitations, highlighted the intricacies of large-scale environmental modeling. Ultimately, our study serves as a stepping stone towards a future where wildfires can be more predictably managed, even as we recognize the importance of continuous model refinement in the face of changing climate dynamics.

References

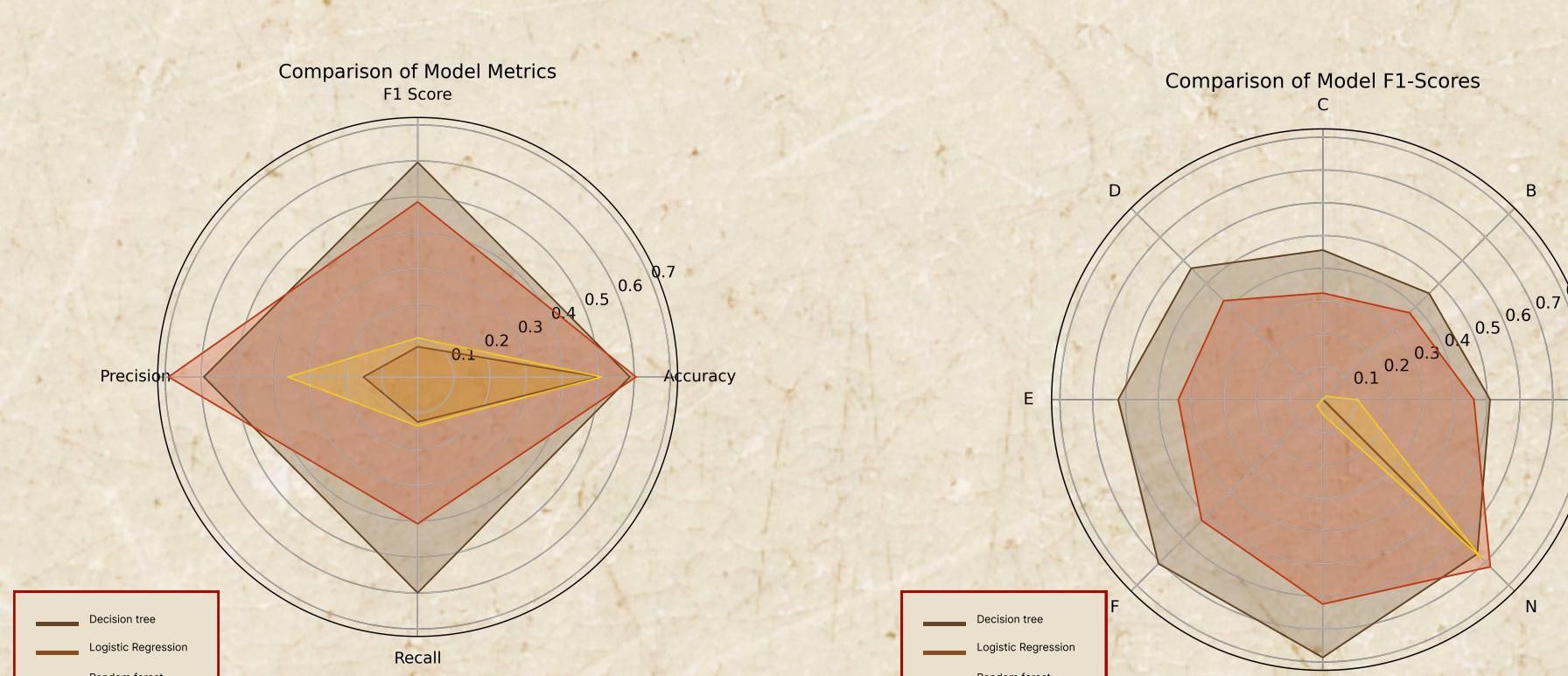
- [1] Short, Koen C. 2020. Spatial wildfire occurrence data for the United States, 1992-2018 [FAO_FOD_20201617]. 5th Edition. Fort Collins, CO: Forest Service Research Data Archive.
- [2] <https://meteostat.net/> (Meteostat)
- [3] Meteostat. 2019. Meteostat API Documentation.
- [4] Montreuil, QC, Canada. 2010. pp. 131-136, doi: 10.1109/AM.2010.5665809.
- [5] Montreuil, QC, Canada. 2010. pp. 479-482, doi: 10.1109/BigData.2012.6278780.
- [6] El-Boraei, R., L. Hu, M. Goyal, T. Sankar, M. Irima and Y. Chen. "Next Day Wildfire Spread: A Machine Learning Dataset to Predict Wildfire Spreading From Remote-Sensing Data," Jan. 2021, Accession Jun. 16, 2023.
- [7] El-Boraei, R., L. Hu, M. Goyal, T. Sankar, M. Irima and Y. Chen. "A New dataset and machine learning approach."Younes Oulad Sayed a, Hajar Mousamri b, Hassan Al Moatassime.

06 RESULTS

Following the inclusion of additional predictors, our models were retrained and evaluated. The tuned Random Forest emerged as the top performer yet again, this time achieving an accuracy of 70%. This substantial improvement underscored the importance of feature engineering and selection in model performance.

In conclusion, our methodology was iterative, beginning with a basic model and feature set, and gradually expanding based on insights and performance metrics. The tuned Random Forest model, when trained with an extended set of predictors, proved to be the most effective in predicting wildfire occurrences and their sizes.

07 DISCUSSION



08 CHALLENGES

The findings of our project presented a series of insightful revelations. Surprisingly, the decision tree emerged as the most effective model. While meteorological factors are undeniably crucial in predicting wildfire occurrences, our analysis underscored the significant influence of geographical location. By integrating data points like latitude, longitude, date, and other geographical attributes, the predictive accuracy experienced a remarkable improvement. This jump in accuracy is a testament to the profound interplay between geographical, temporal, and meteorological factors in the behavior of wildfires. Seasonal variations were also clearly manifested in our results. For instance, fire occurrences dwindled during winter, aligning with the intuitive understanding of the inhibitive role cold weather plays in fire spread. From a wildfire protection and management perspective, the implications are manifold. Our findings suggest that merely relying on meteorological predictors might not furnish a holistic picture. Incorporating geographical data can aid in devising more localized and effective fire prevention strategies. Moreover, understanding the temporal patterns, such as reduced fires in winter, can help authorities optimize resource allocation during peak fire seasons. This could lead to better-preparedness and more efficient response mechanisms, ultimately safeguarding ecosystems and communities alike.

Large Dataset Size: Navigating a dataset of 2.5 million rows required strategic preprocessing, especially given the volume of missing values in columns like 'tmax', 'tmin', and 'prec'.