

Canada Wildfire Analysis and Prediction: Project Report

Problem Definition

Wildfires have a profound impact on Canada's ecosystems, economy, and public safety, with climate change driving an increase in their frequency and severity. By analyzing wildfire data from NASA's Fire Information for Resource Management System (FIRMS), which provides insights into fire brightness and locations, and combining it with historical weather data from Open-Meteo, we aim to explore patterns and build predictive models. These models will help understand and forecast wildfire occurrences based on factors such as location, weather conditions, and time.

The datasets from FIRMS and Open-Meteo are substantial, often exceeding the capacity of a single processor. To effectively analyze this big data, we will leverage distributed computing using Python and Apache Spark, enabling scalable and efficient data processing.

Integrating satellite data from NASA FIRMS with weather data from Open-Meteo posed significant challenges. The datasets differed in temporal and spatial resolutions, requiring meticulous alignment and preprocessing to ensure accuracy. Temporal variations, such as mismatched timestamps, and spatial discrepancies, like differing grid sizes, had to be resolved to create a unified dataset suitable for analysis. Another challenge was understanding the multivariate dependencies between environmental factors such as dryness, temperature, and precipitation. These variables interact in complex ways, and capturing their interdependencies required advanced analysis techniques. Identifying meaningful patterns while accounting for the influence of multiple factors added further complexity. Finally, optimizing predictive models was critical when working with high-dimensional and noisy data. Selecting algorithms that could balance computational efficiency and accuracy, while minimizing the risk of overfitting, demanded careful experimentation and tuning. Achieving robust predictions required iterative testing and validation to ensure generalization to unseen data.

Methodology

Data Sources:

For access to historical data, Open Meteo's api was used. Instead of focusing on hundreds of cities across Canada, our team decided on the cities that usually have the most wildfires. The cities which we chose were: Fort McMurray, Kelowna, Kamloops, Prince George, Vancouver Island, Lytton, Penticton, Williams Lake, Grande Prairie, and Edson. The dataset

included features such as date, city, latitude, longitude, variations of temperature, precipitation types and amounts, and wind speed.

In order to access historical wildfire data, we sourced it from NASA FIRMS (MODIS). This dataset includes wildfire-specific data such as brightness, frp, confidence, and bright_t31. After requesting data for Canada, we received the data and joined it with weather data from Open Meteo. This join occurred on longitude, latitude, and date. This join resulted in being able to view weather statistics from days with an active wildfire and non-active wildfires.

Data Prep:

While we have significant and important data for wildfire analysis, adding more meaningful features would result in a more thorough analysis. There are many factors which play a big part in wildfire occurrence, such as soil moisture. When soil moisture is low, we can assume that the vegetation is also dry, which leads to stronger wildfires as dry vegetation acts as a fuel for fire. If the value for soil moisture is higher, we can make estimations for precipitation levels, seasonal changes, and temperature approximations.

Another important feature for wildfire analysis is the intensity of a wildfire. By combining MODIS data: brightness, confidence, and frp, the result yields intensity. Brightness tells us how intense the light is coming from the wildfire. Confidence delivers how certain the satellite is for the occurrence of the wildfire. FRP is a measure of how much heat is radiating from the fire itself. Then, the result, intensity, gives us the combined assumption about the fire itself. A higher intensity value means the fire is burning hotter and producing very high heat. A lower value could possibly reflect a weak or dying wildfire.

By adding all the varieties of precipitation and their respective amounts, we designed a new feature called cumulative precipitation. This assists in many aspects of a wildfire such as intensity and soil moisture. Higher levels of precipitation often dictate lower chances of a wildfire happening. Lower levels contribute to low soil moisture, which in return can result in more wildfire occurrences and higher intensity fires.

While the summary above does not include all of the new features which were added during feature engineering, they dictate the thought behind our analysis. Other features which were not mentioned include: humidity, dryness, risk, and precipitation-sun ratio. All these features were constructed for more meaningful analysis and each plays its importance in wildfire occurrence.

Data Analysis:

Our data analysis was split into 3 sections: feature analysis, analysis per city, and yearly trends. Feature analysis focused on which features correlated to more intense fires or a

wildfire event. During this analysis, finding out which feature was the best for wildfire prediction was also calculated. Majority of this analysis occurred in Spark using dataframes aggregations. Plots were constructed using matplotlib and seaborn. Each data frame was transformed into a pandas dataframe before plotting. Another aspect which was mentioned before was feature importance. Our team conducted 2 feature selections using PCA and a RandomForestRegressor.

Analysis by city was also one of our main focuses. Our approach for this analysis was to compare cities and see how their wildfires differed. While wildfire occurrence was our main topic for this, we also wanted to compare a variety of features and their importance to each city's fire risk. Another important aspect to our analysis was yearly trends. By viewing the different changes over the years per each city, we can visualize how the overall changes in climate contribute to either stronger or weaker wildfire seasons.

Machine Learning:

Machine learning was also split into 3 sections: Regression, Classification, and Neural Networks. For classification, both RandomForestClassifier and GBTCClassifier were used in a Spark setting. Before any optimization, both models held an accuracy of 87% for predicting a wildfire event. After hyperparameter tuning for both models, they both yield close to 95% in wildfire predictions. The hyperparameter tuning for the GBTCClassifier includes learning rate alterations and for RandomForestClassifier: changing values for the number of trees and the maximum size of each tree.

The neural network regressor was exceptionally effective, achieving an impressive coefficient of determination (R^2) of 0.991 and a Root Mean Squared Error (RMSE) of 15.382, highlighting its precision in predicting brightness levels. XGBoost also performed robustly with an RMSE of 16.14 and an R^2 of 0.99, demonstrating high predictive accuracy. In contrast, LightGBM showed decent performance, though slightly less effective, with an RMSE of 15.61 and an R^2 of 0.837. Other regression models, including Linear Regression and Random Forest, did not perform as well, reflecting higher error values in their predictions. Additionally, our neural network classifier, a Multilayer Perceptron (MLP), after hyperparameter tuning, reached an accuracy of 0.999, showing excellent performance without any signs of overfitting, underscoring the strength of our predictive modeling approach.

Repository:

Our GitHub repository includes our entire data pipeline. This includes scripts for feature engineering and the resulting dataset from the merging of all relevant data. It also holds all analysis code including their respective visualizations. Each machine learning model is also included for further predictions and tuning.

Problems Encountered

At the outset of the project, identifying reliable and relevant data sources was a significant challenge. Early attempts with various datasets revealed issues such as incorrect or missing values, limited support for large downloads, and insufficient documentation. These shortcomings hindered initial progress, requiring a shift toward well-documented and robust datasets like those from NASA FIRMS and Open-Meteo, which aligned better with the project goals.

Data quality also posed hurdles, as missing and noisy data disrupted early modeling efforts. Weather data gaps and outliers in fire brightness values had to be addressed to improve data integrity. Interpolation techniques were implemented for missing weather values, while outlier detection and removal ensured the reliability of brightness data. These preprocessing steps enhanced the foundation for analysis and modeling.

Additionally, not all derived features proved useful for prediction. For example, variables like wind gust showed minimal contribution to model performance. To address this, feature importance rankings were employed to identify the most relevant inputs, and dimensionality reduction techniques streamlined the feature set. Computational challenges further compounded these issues, particularly when training neural networks, which required significant resources. Leveraging distributed environments like PySpark and optimizing batch processing allowed for efficient model training and scalability.

Results

Our analysis revealed that brightness is a strong predictor of wildfire occurrence, particularly under dry and warm conditions. High brightness values, particularly under dry and warm conditions, strongly correlated with extreme weather patterns, indicating that brightness levels can serve as a valuable indicator for fire risk, especially in regions experiencing heightened temperatures and dry spells.

Geographically, cities such as Kamloops, Penticton, and Kelowna exhibited consistently high brightness values and elevated wildfire risk. These findings emphasize the need for targeted interventions and preventative measures in these areas, where environmental and climatic conditions exacerbate the potential for wildfire outbreaks. The combination of high brightness, persistent dryness, and temperature spikes makes these cities particularly vulnerable to wildfires.

From an environmental perspective, increasing temperatures and decreasing precipitation were found to directly influence wildfire patterns. The rising temperatures coupled with

reduced rainfall are accelerating the frequency and intensity of wildfires, underscoring the importance of understanding these environmental trends for more effective wildfire management strategies.

For our machine learning model performance, we employed a variety of machine learning models to predict wildfire occurrences based on brightness data. Neural networks proved the most effective, achieving a Root Mean Squared Error (RMSE) of 15.38 and a coefficient of determination (R^2) of 0.991 for brightness prediction. This outperformed other methods in handling the dataset's complexity, offering the most reliable predictions.

For classification tasks, the neural network (multilayer perceptron) has the highest accuracy of 99.1 % followed by the Random Forest model demonstrated the highest accuracy at 94.9%, outperforming logistic regression and support vector machines (SVM), which performed well but were slightly less precise. In regression tasks, the best performance was multilayer perceptron (neural networks) with r^2 score of 0.99 and rmse score of 15.38, followed by XGBoost showed strong performance with an RMSE of 16.14 and an R^2 of 0.99 for brightness prediction, while LightGBM delivered an RMSE of 15.61 and an R^2 of 0.837, making it a good alternative but not quite as effective as XGBoost.

Project Summary

★ Getting the Data (3 points)

Data was sourced from NASA FIRMS and Open Meteo APIs. We submitted requests to NASA FIRMS for Canadian data and applied logic to extract data for the top ten cities based on latitude and longitude values. APIs were utilized to obtain weather data using the libraries `openmeteo-requests`, `requests-cache`, `retry-requests`, `numpy`, and `pandas`. We acquired detailed satellite and weather data, including brightness and climate variables.

★ ETL: Extract-Transform-Load (2 points)

We performed extensive cleaning and merging of datasets, addressing missing values and aligning spatial-temporal dimensions. New features such as dryness index, precipitation ratios, and temperature range were engineered to enhance model inputs.

★ Problem Definition (3 points)

We clearly defined wildfire prediction as the problem, focusing on brightness as the target metric and the classification of wildfire occurrences (whether a fire happens or not). This work was motivated by the growing need for wildfire risk management in the context of climate change.

★ Algorithmic Work (4 points)

We implemented the following machine learning models for regression to predict brightness (fire intensity) in MODIS data:

- Random Forest, Linear Regression , XGBoost, Multilayer Perceptron

For binary classification of the column `in_modis` (indicating whether a fire occurs), we applied:

- Random Forest Classifier, Gradient-Boosted Tree Classifier, Logistic Regression , Multilayer Perceptron

Machine learning algorithms were implemented with an emphasis on neural networks for brightness prediction and classification.

We tuned models and validated their performance using metrics such as RMSE and accuracy.

★ **Bigness/Parallelization (3 points)**

We leveraged PySpark for distributed data processing and model training. Workflows were optimized for scalability to handle larger datasets.

★ **UI: User Interface (0 points)**

This project did not include a user interface or interactive front end.

★ **Visualization (3 points)**

We developed insightful visualizations to analyze trends in temperature, precipitation, and wildfire activity. Heat maps and temporal graphs were used to illustrate correlations between features and wildfire risk.

An attached report provides an in-depth explanation of our analysis.

★ **Technologies (2 point)**

We learned and applied PySpark ML for distributed machine learning workflows, employing algorithms such as SVM, logistic regression, linear regression, and neural network perceptron. Additionally, we gained experience with feature engineering techniques, including PCA and PySpark tree importances, as well as advanced regression methods.

Total: 20 points

This summary captures our project's emphasis on ETL, algorithmic development, and visual analysis, with moderate attention to data acquisition and scalability.

While we did not focus on UI development, the strengths in data processing, modeling, and visualization form the core contributions of our project