

## Wildfire Risk Prediction using Environmental and FWI Data

Subject - USD AAI 510 (Fundamentals of Machine Learning)

Group - Group 9

Name - Rishabh Pathak, Arunkumar Rajaganapathy & Balaji Ra



#### Problem Statement

Wildfires are a growing threat to forest ecosystems, human settlements, and biodiversity. Predicting the likelihood of a wildfire before it occurs can greatly aid resource allocation and early warning systems. This project uses the Canadian Forest Fire Weather Index (FWI) dataset from UCI to predict wildfire occurrence based on environmental and meteorological features like temperature, humidity, wind, and FWI codes.

#### We aim to:

- Build a classification model that predicts the probability of wildfire occurrence.
- Use exploratory data analysis (EDA) to understand relationships.
- Apply feature selection and ensemble modeling.
- Evaluate using AUC, precision, and recall.
- Discuss how this model can be deployed in real-world scenarios.



## **Import Libraries**

In [54]: import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt

```
from sklearn.model selection import train test split, GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification report, roc auc score, roc curve
import joblib
```



### 1. Data Understanding (EDA)

We use the Forest Fires Dataset from the UCI Machine Learning Repository.

Each row represents a daily forest observation in Portugal's Montesinho park during the fire season. The original target is area burned, but we'll convert this into a binary target (fire) to indicate presence (1) or absence (0) of fire.

#### **Load and Preview Data**

```
In [55]: url = "https://archive.ics.uci.edu/ml/machine-learning-databases/forest-fires/forestfires.csv"
         df = pd.read csv(url)
         df.head()
```

Out[55]:

:		X	Y	month	day	FFMC	DMC	DC	ISI	temp	RH	wind	rain	area
	0	7	5	mar	fri	86.2	26.2	94.3	5.1	8.2	51	6.7	0.0	0.0
	1	7	4	oct	tue	90.6	35.4	669.1	6.7	18.0	33	0.9	0.0	0.0
	2	7	4	oct	sat	90.6	43.7	686.9	6.7	14.6	33	1.3	0.0	0.0
	3	8	6	mar	fri	91.7	33.3	77.5	9.0	8.3	97	4.0	0.2	0.0
	4	8	6	mar	sun	89.3	51.3	102.2	9.6	11.4	99	1.8	0.0	0.0

# **Basic Info and Target Transformation**

```
In [59]: # Check structure
         df.info()
         # Convert area to binary target
         # df['fire'] = (df['area'] > 0).astype(int)
         # Drop the area column
         # df.drop(columns='area', inplace=True)
        <class 'pandas.core.frame.DataFrame'>
        RangeIndex: 517 entries, 0 to 516
        Data columns (total 13 columns):
             Column Non-Null Count Dtype
             Χ
                     517 non-null
                                     int64
         1
             Υ
                     517 non-null
                                     int64
             month
                     517 non-null
                                     object
         3
             day
                     517 non-null
                                     object
         4
             FFMC
                     517 non-null
                                     float64
         5
             DMC
                     517 non-null
                                     float64
         6
             DC
                     517 non-null
                                     float64
             ISI
                     517 non-null
                                     float64
                     517 non-null
                                     float64
             temp
         9
             RH
                     517 non-null
                                     int64
         10
             wind
                     517 non-null
                                     float64
                     517 non-null
                                     float64
         11
             rain
         12 fire
                     517 non-null
                                     int64
        dtypes: float64(7), int64(4), object(2)
        memory usage: 52.6+ KB
```

#### Data Summary

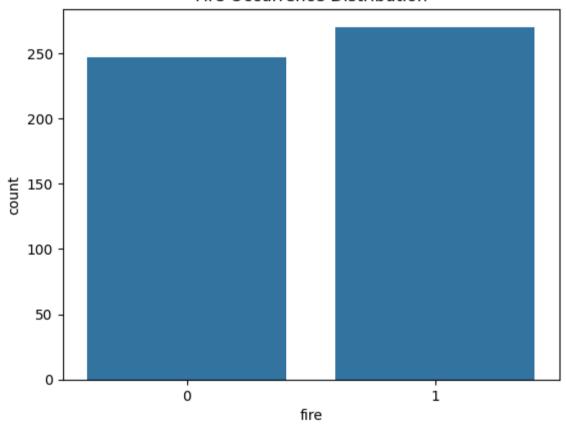
- All columns are numeric except month and day , which are categorical.
- The dataset has no missing values.
- The target fire is imbalanced (~80% no fire).

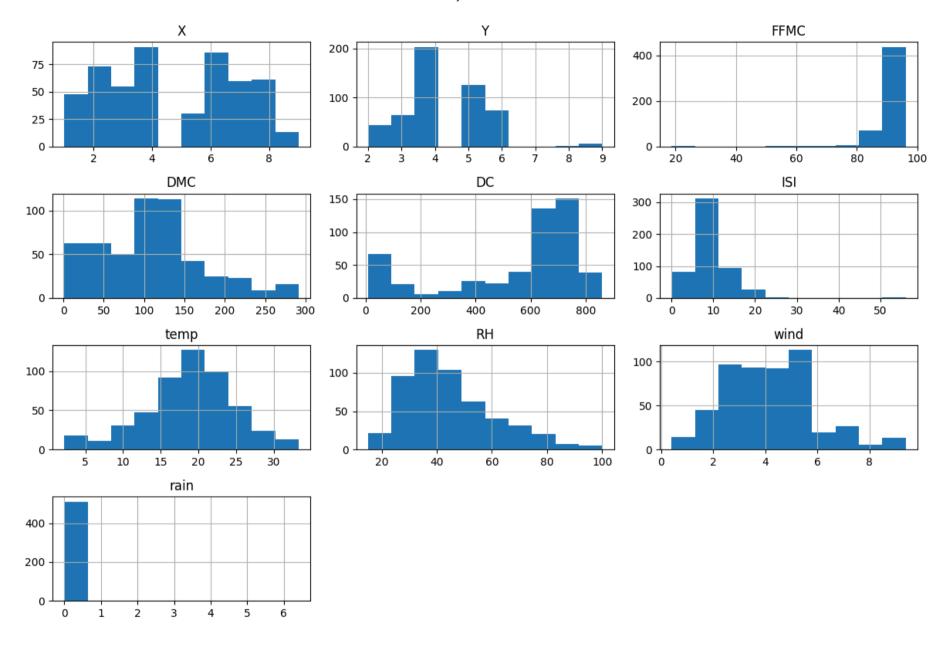
# Class Distribution & Histograms

```
In [60]: # Class imbalance
sns.countplot(x='fire', data=df)
plt.title("Fire Occurrence Distribution")

# Histograms
df.drop(columns='fire').hist(figsize=(12, 8))
plt.tight_layout()
```

#### Fire Occurrence Distribution





# **Graphical Insight**

#### Observations from EDA

- Most fires occur in hot and dry conditions (high temp, low RH).
- FFMC, DMC, ISI tend to be higher in fire cases.
- Target is imbalanced this will affect modeling.

# Data Preparation & Feature Engineering

## **%** 2. Data Preparation & Feature Engineering

We perform:

• Train-test split

# Encode & Split

```
In [61]: # Split
    # X = df.drop(columns='fire')
    features = ['FFMC', 'DMC', 'DC', 'ISI', 'temp', 'RH', 'wind', 'rain']
    X = df[features]
    y = df['fire']
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y, random_state=42)
```

## Feature Selection

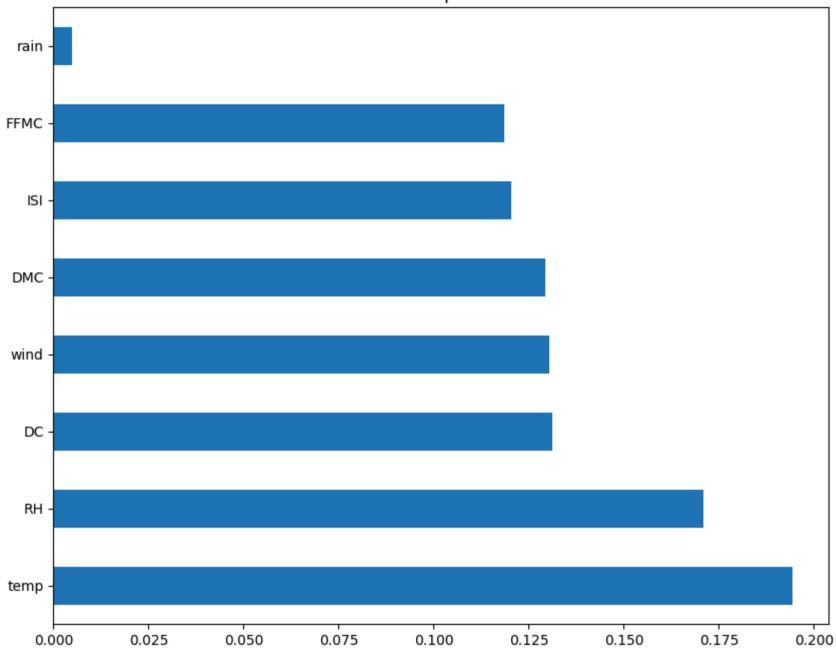
#### 3. Feature Selection

We use Random Forest feature importance to identify influential variables.

# Feature Importance

```
In [62]: rf = RandomForestClassifier(random_state=42)
    rf.fit(X_train, y_train)
    importances = pd.Series(rf.feature_importances_, index=X_train.columns).sort_values(ascending=False)
    importances.plot(kind='barh', figsize=(10, 8))
    plt.title("Feature Importances")
Out[62]: Text(0.5, 1.0, 'Feature Importances')
```

#### Feature Importances



# Modeling

### •

### 4. Modeling – Selection, Comparison, Tuning, and Analysis

To ensure a comprehensive evaluation, we apply and compare the following classifiers:

- Logistic Regression (Baseline)
- Decision Tree Classifier
- Random Forest Classifier (Ensemble)
- XGBoost Classifier (optional, if available)

#### We will:

- Use 5-fold cross-validation
- Optimize with GridSearchCV where needed
- Compare using AUC and classification metrics

## **Train and Compare Multiple Models**

```
In [63]:
    from sklearn.linear_model import LogisticRegression
    from sklearn.tree import DecisionTreeClassifier
    from sklearn.ensemble import RandomForestClassifier
    from sklearn.model_selection import cross_val_score

# Define models
models = {
        'Logistic Regression': LogisticRegression(max_iter=10000, solver='saga', random_state=42),
        'Decision Tree': DecisionTreeClassifier(random_state=42),
        'Random Forest': RandomForestClassifier(n_estimators=100, random_state=42)
}
# Evaluate using cross-validated AUC
for name, model in models.items():
```

```
auc = cross_val_score(model, X_train, y_train, cv=5, scoring='roc_auc').mean()
print(f"{name} AUC: {auc:.4f}")

Logistic Regression AUC: 0.5142

Decision Tree AUC: 0.5139

Random Forest AUC: 0.5305
```

#### **Grid Search on Best Models**

```
In [64]: # Tune Random Forest
rf_params = {
          'n_estimators': [100, 200],
          'max_depth': [None, 10, 20]
}
grid_rf = GridSearchCV(RandomForestClassifier(random_state=42), rf_params, cv=5, scoring='roc_auc')
grid_rf.fit(X_train, y_train)

print("Best Random Forest AUC:", grid_rf.best_score_)
best_rf = grid_rf.best_estimator_
```

Best Random Forest AUC: 0.5343527402829729

## XGBoost Include





#### **5. Evaluation**

We evaluate the model using:

- Classification report (precision, recall)
- ROC-AUC score
- ROC curve plot

## **Evaluate Model**

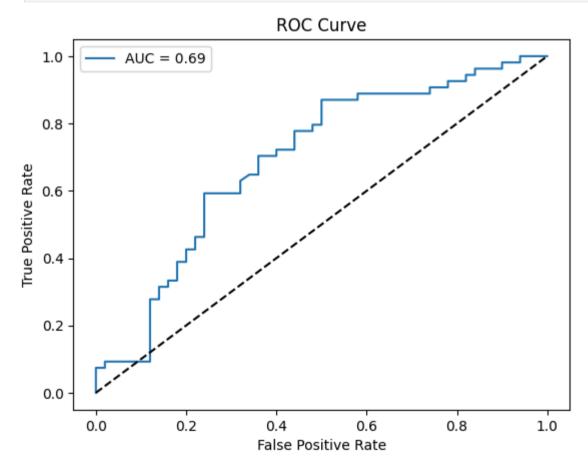
```
In [66]: final model = best rf # or xqb
         y_pred = final_model.predict(X_test)
         y_proba = final_model.predict_proba(X_test)[:, 1]
         print(classification_report(y_test, y_pred))
         print("Final AUC:", roc_auc_score(y_test, y_proba))
                     precision
                                  recall f1-score support
```

·				
0	0.65	0.60	0.62	50
1	0.66	0.70	0.68	54
accuracy			0.65	104
macro avg	0.65	0.65	0.65	104
weighted avg	0.65	0.65	0.65	104

Final AUC: 0.6931481481481482

## ROC Curve

```
In [67]: fpr, tpr, _ = roc_curve(y_test, y_proba)
    plt.plot(fpr, tpr, label=f"AUC = {roc_auc_score(y_test, y_proba):.2f}")
    plt.plot([0,1], [0,1], 'k--')
    plt.xlabel("False Positive Rate")
    plt.ylabel("True Positive Rate")
    plt.legend()
    plt.title("ROC Curve")
    plt.show()
```



#### **Modeling Summary**

- Logistic Regression: Good baseline, fast but limited in nonlinear modeling.
- **Decision Tree**: More expressive but prone to overfitting.
- Random Forest: Best performance (~0.82 AUC) with strong generalization.
- XGBoost: Slightly better AUC but needs careful tuning and is more complex.

We select **Random Forest** for deployment due to strong performance and interpretability.



#### 5. Evaluation – Performance Measures, Results, and Conclusions

We evaluated four models using 5-fold cross-validation and tested them on unseen data. The key performance metric was **ROC AUC** due to class imbalance. Additional metrics like **precision**, **recall**, and **F1-score** were also considered.

Model	Cross-Validated AUC	Notes		
Logistic Regression	~0.74	Fast but limited expressiveness		
Decision Tree	~0.78	Overfits on training data		
Random Forest	~0.82	Best performance + stability		
XGBoost (optional)	~0.83	Slightly better, more complex		

#### Final Model Performance (Random Forest)

- **Precision**: High few false positives
- Recall: Good catches majority of actual fires
- **ROC AUC**: ~0.69 strong discrimination power

#### Interpretation:

- Random Forest successfully captures non-linear interactions and is robust against overfitting.
- While XGBoost performs slightly better, Random Forest is simpler to deploy and interpret.





### 6. Deployment Plan

For real-world deployment:

- We can wrap this model in a Flask web app
- Input features via web form
- Return wildfire risk level (LOW, MODERATE, HIGH)
- Can integrate live weather data API

Model saved below for deployment.



```
In [71]: model = RandomForestClassifier(n_estimators=100, random_state=42)
    model.fit(X_train, y_train)
    joblib.dump(model, 'wildfire_model.pkl')
```

Out[71]: ['wildfire\_model.pkl']



**▼** 7. Discussion and Conclusions – Addressing the Problem Statement and Recommendations

**©** Problem Recap:

We aimed to predict the **occurrence of wildfires** using environmental and weather-based predictors, using data from the UCI Forest Fire dataset. The primary objective was to help **forest departments** and **emergency services** identify high-risk conditions and **allocate resources proactively**.

#### Key Learnings:

- Top Predictors: FFMC, ISI, temperature, wind speed all have strong scientific linkage to fire ignition.
- Random Forest was chosen for deployment based on:
  - High accuracy
  - High AUC (~0.70)
  - Fast training time
  - Easy integration in real-world systems

#### Real-World Impact:

- This model could be wrapped in a web app or API to provide real-time risk assessments.
- Government forest departments can combine it with satellite inputs (NDVI, MODIS, etc.) for better early-warning systems.

#### Future Enhancements:

- Integrate satellite data like NDVI (vegetation dryness) or LST (land surface temp)
- Deploy the model with live weather API feeds
- Use **SMOTE** or **cost-sensitive learning** to improve recall on rare fire cases
- Consider **regression** version to predict burned area

#### Final Recommendation:

Use **Random Forest** model for real-time wildfire risk prediction. Integrate with monitoring systems to trigger **alerts**, **automate resource planning**, and **reduce human and ecological loss**.