

# Analisis Dataset Forest Fires - Penjelasan Tujuan dan Algoritma

## Bagian 1: Penjelasan Tujuan dan Algoritma

### A. Tujuan Proyek

#### 1. Tujuan Utama

Tujuan utama dari proyek analisis dataset Forest Fires ini adalah untuk **membangun model prediktif yang dapat memprediksi luas area hutan yang terbakar (burned area)** berdasarkan kondisi meteorologi dan faktor spasial-temporal. Ini merupakan masalah **regresi** karena variabel target 'area' merupakan nilai kontinyu yang mewakili luas area terbakar dalam hektar (0.00 - 1090.84 ha).

#### 2. Tujuan Spesifik

- **Prediksi Dini:** Mengembangkan sistem peringatan dini untuk memperkirakan potensi kerusakan kebakaran hutan berdasarkan kondisi cuaca
- **Manajemen Risiko:** Membantu pihak berwenang dalam mengalokasikan sumber daya pemadaman kebakaran secara efisien
- **Pemahaman Pola:** Mengidentifikasi faktor-faktor meteorologi yang paling berpengaruh terhadap intensitas kebakaran hutan
- **Optimasi Respons:** Memberikan informasi untuk perencanaan strategi pencegahan dan penanggulangan kebakaran

#### 3. Manfaat Praktis

- Mengurangi kerugian ekonomi dan ekologi akibat kebakaran hutan
- Meningkatkan efektivitas alokasi tim pemadam kebakaran
- Mendukung kebijakan konservasi hutan yang berbasis data

### B. Algoritma yang Digunakan

Mengingat karakteristik dataset dan nature dari masalah regresi ini, saya merekomendasikan penggunaan **beberapa algoritma machine learning** dengan pendekatan komparatif:

#### 1. Random Forest Regressor (Algoritma Utama)

##### Alasan Pemilihan:

- **Robust terhadap outliers:** Dataset ini memiliki distribusi target yang sangat skewed (banyak nilai 0 dan beberapa nilai ekstrem tinggi)

- **Feature importance:** Dapat mengidentifikasi variabel meteorologi mana yang paling berpengaruh
- **Non-linear relationships:** Mampu menangkap hubungan kompleks antara kondisi cuaca dan intensitas kebakaran
- **Overfitting prevention:** Built-in regularization melalui ensemble learning

## 2. Gradient Boosting (XGBoost) (Algoritma Pembanding)

### Alasan Pemilihan:

- **Performance tinggi:** Umumnya memberikan akurasi prediksi yang excellent
- **Handling missing values:** Robust terhadap data yang tidak lengkap
- **Feature selection otomatis:** Secara otomatis memberikan bobot pada fitur yang penting
- **Scalability:** Efisien untuk dataset berukuran sedang seperti ini (517 instances)

## 3. Support Vector Regression (SVR) (Algoritma Pembanding)

### Alasan Pemilihan:

- **Reference benchmark:** Paper asli dataset ini menggunakan SVM dan mendapatkan hasil terbaik (MAD: 12.71)
- **Kernel flexibility:** Dapat menangkap pola non-linear melalui RBF kernel
- **Robust to outliers:** Menggunakan epsilon-insensitive loss function
- **Regularization:** Built-in regularization untuk mencegah overfitting

## C. Justifikasi Algoritma

### 1. Mengapa Ensemble Methods (Random Forest & XGBoost)?

- Dataset memiliki **high variance** dalam target variable (area: 0 - 1090.84)
- **Multiple weak predictors** dapat dikombinasikan untuk membentuk strong predictor
- **Bias-variance trade-off** yang optimal untuk masalah regresi kompleks
- **Feature interaction** antara variabel meteorologi dapat ditangkap dengan baik

### 2. Mengapa SVR sebagai Baseline?

- **Historical benchmark:** Paper asli menggunakan SVM dengan hasil terbaik
- **Mathematical foundation:** Solid theoretical background untuk regression tasks
- **Comparison purpose:** Memberikan baseline performance untuk membandingkan algoritma modern

### 3. Pertimbangan Khusus untuk Dataset Ini:

- **Skewed distribution:** Target variable sangat skewed ke 0, sehingga diperlukan transformasi logaritmik  $\ln(x+1)$
- **Weather dependency:** Kondisi meteorologi memiliki hubungan non-linear yang kompleks
- **Temporal patterns:** Variabel month dan day mengindikasikan seasonal patterns
- **Spatial correlation:** Koordinat X,Y menunjukkan spatial dependency

#### 4. Evaluation Strategy:

Menggunakan **multiple metrics** untuk evaluasi yang komprehensif:

- **MAD (Mean Absolute Deviation):** Fokus pada error absolut, sesuai dengan paper original
- **RMSE (Root Mean Square Error):** Memberikan penalty lebih besar untuk error besar
- **R<sup>2</sup> Score:** Mengukur proportion of variance yang dapat dijelaskan model
- **Cross-validation (10-fold):** Sesuai metodologi paper asli untuk hasil yang robust

#### 5. Preprocessing Requirements:

- **Log transformation:**  $\ln(\text{area} + 1)$  untuk mengatasi skewed distribution
- **Feature scaling:** StandardScaler untuk SVR, tidak diperlukan untuk tree-based methods
- **Categorical encoding:** One-hot encoding untuk month dan day
- **Outlier detection:** Identifikasi dan handling extreme values

### Kesimpulan Strategis

Pendekatan multi-algoritma ini memberikan beberapa keuntungan:

1. **Comprehensive analysis:** Setiap algoritma memberikan insight yang berbeda
2. **Performance comparison:** Dapat menentukan algoritma terbaik untuk kasus spesifik ini
3. **Robust prediction:** Ensemble dari multiple models dapat memberikan prediksi yang lebih stabil
4. **Scientific rigor:** Mengikuti metodologi yang established dalam literatur machine learning

Dengan kombinasi ketiga algoritma ini, diharapkan dapat menghasilkan model prediktif yang akurat dan reliable untuk mendukung manajemen kebakaran hutan yang lebih efektif.