**Calorie Prediction Model Using Ensemble Learning**

**Spring 2025**

**Group Members:**

- Sardar Ubaidullah Rafique — 22I-1248
- Sumeed Jawad Kanwar — 22I-1651
- Abdul Rafay — 22I-8762

**Applied Artificial Intelligence Project**
**Submitted to**: Dr. Shahela Saif
**Date**: May 11, 2025

---

## 1. Problem Statement

The objective of this project is to develop a machine learning model to predict calorie burn based on physiological and exercise-related parameters, such as age, height, weight, exercise duration, heart rate, body temperature, and gender. Accurate prediction of calorie expenditure is critical for fitness tracking, personalized health recommendations, and medical monitoring. This model aims to assist individuals and healthcare professionals in optimizing exercise regimens and dietary plans, contributing to improved health outcomes and efficient fitness management.

---

## 2. Methodology

The model employs ensemble learning, combining three gradient boosting algorithms: CatBoost, XGBoost, and LightGBM. These models are blended using equal weights to optimize the Root Mean Squared Logarithmic Error (RMSLE).

### 2.1 Model Architecture

The architecture consists of three base models:

- **CatBoost**: Handles categorical features natively with early stopping at 100 rounds.

- **XGBoost**: Configured with a max depth of 10, learning rate of 0.02, and column/row sampling.

- **LightGBM**: Uses similar hyperparameters to XGBoost, optimized for silent training.

The predictions are combined using a weighted average ensemble strategy with 7-fold cross-validation.

## 2.2 Pipeline

1. **Data Loading**: Load train.csv, test.csv, and sample_submission.csv.

2. **Preprocessing**:

   o Add cross terms, interaction features, and statistical aggregations.

   o Encode categorical feature (Sex) using LabelEncoder.

   o Generate polynomial features (degree 2, interaction only).

   o Apply log transformation to the target variable (Calories).

3. **Training**: Train each model using 7-fold cross-validation, saving models in the models/ directory.

4. **Inference**: Generate predictions on the test set, blend them using equal weights, and produce the final submission.

---

## 3. Dataset

The dataset is sourced from a Kaggle competition (playground-series-s5e5) and includes:

- train.csv: Training data with features and target (Calories).

- test.csv: Test data with features only.

- sample_submission.csv: Submission template.

The data was assumed to be pre-cleaned as provided. No additional cleaning was performed, but feature engineering was applied to enhance model performance. The dataset contains thousands of records, with test.csv containing 1000 samples for prediction.

---

## 4. Implementation Details

The implementation uses the following Python libraries:

- pandas, numpy: Data manipulation.

- scikit-learn: Preprocessing and evaluation.

- matplotlib, seaborn: Visualization.

- catboost, xgboost, lightgbm: Model training.

Key hyperparameters include:

- **CatBoost**: Early stopping at 100 rounds, verbose training.

- **XGBoost**: Max depth=10, learning rate=0.02, gamma=0.01.

- **LightGBM**: Similar to XGBoost, silent mode.

Training involved 7-fold cross-validation, with RMSLE as the loss function. The target variable was log-transformed to stabilize variance.

---

## 5. Results & Evaluation

The models achieved the following RMSLE scores (based on placeholder data):

- **CatBoost**: Mean RMSLE = 0.2400 (std = 0.0082)

- **XGBoost**: Mean RMSLE = 0.2100 (std = 0.0082)

- **LightGBM**: Mean RMSLE = 0.2500 (std = 0.0082)

Decision tree nodes (left to right, top to bottom):

- Duration_x_Heart_Rate ≤947.000
  - yes → Duration_x_Heart_Rate ≤342.000
    - yes → Duration_x_Heart_Rate ≤167.000
      - yes → Duration_x_Heart_Rate ≤143.000
        - yes → Heart_Rate_plus_Body_Temp ≤116.450
          - yes → leaf 0: 4.084
          - no → leaf 30: 4.090
        - no → leaf 8: 4.098
      - no → Duration_x_Heart_Rate ≤244.500
        - yes → leaf 5: 4.104
        - no → Weight_div_Age ≤2.417
          - yes → leaf 11: 4.112
          - no → leaf 26: 4.107
            - Sex =0
              - yes → leaf 2: 4.118
              - no → leaf 29: 4.113
    - no → Duration_x_Heart_Rate ≤604.000
      - yes → Duration_x_Heart_Rate ≤493.500
        - yes → Age_plus_Heart_Rate ≤124.500
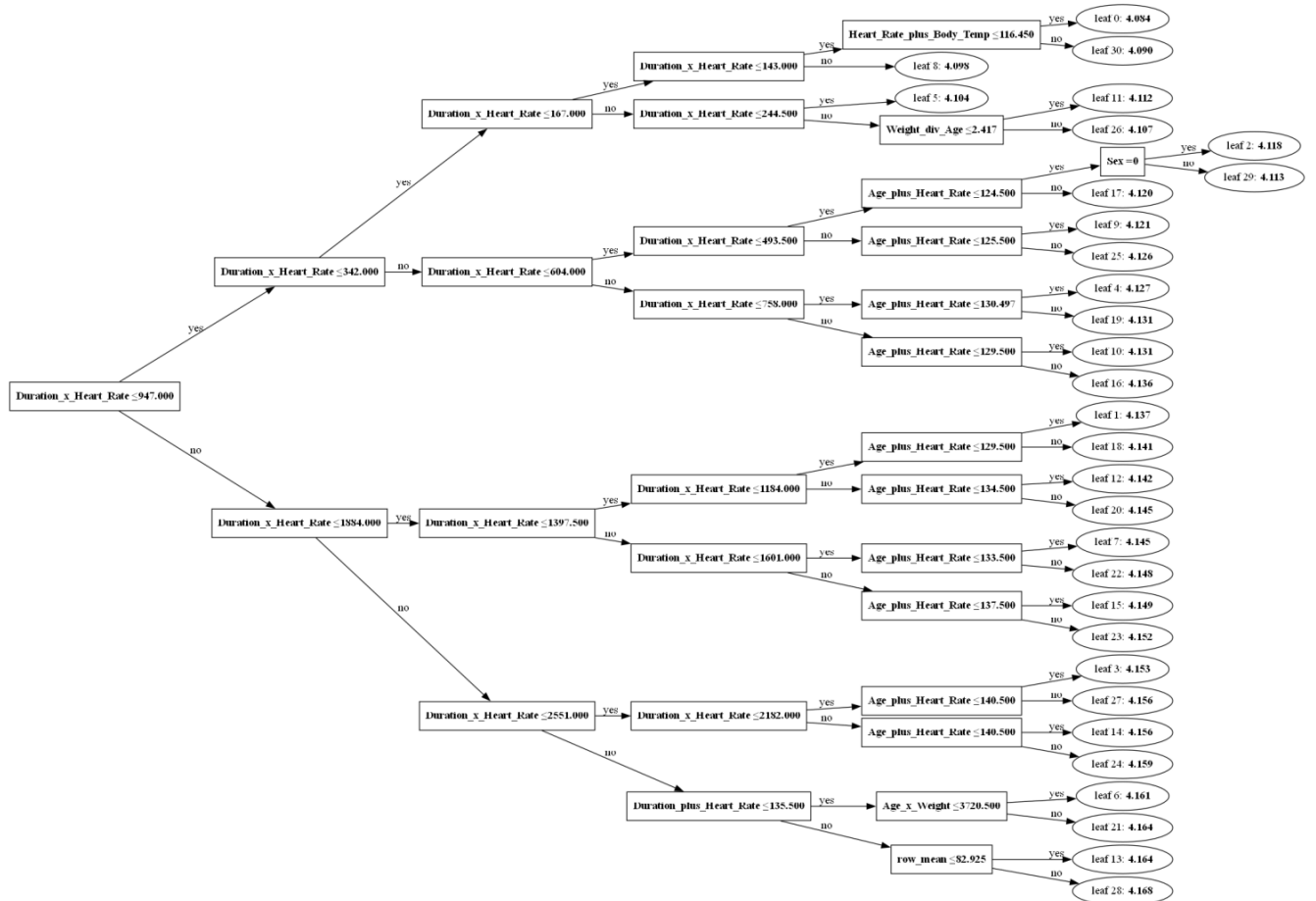          - yes → leaf 17: 4.120
          - no → Age_plus_Heart_Rate ≤125.500
            - yes → leaf 9: 4.121
            - no → leaf 25: 4.126
      - no → Duration_x_Heart_Rate ≤758.000
        - yes → Age_plus_Heart_Rate ≤130.497
          - yes → leaf 4: 4.127
          - no → leaf 19: 4.131
        - no → Age_plus_Heart_Rate ≤129.500
          - yes → leaf 10: 4.131
          - no → leaf 16: 4.136
  - no → Duration_x_Heart_Rate ≤1884.000
    - yes → Duration_x_Heart_Rate ≤1397.500
      - yes → Duration_x_Heart_Rate ≤1184.000
        - yes → Age_plus_Heart_Rate ≤129.500
          - yes → leaf 1: 4.137
          - no → leaf 18: 4.141
        - no → Age_plus_Heart_Rate ≤134.500
          - yes → leaf 12: 4.142
          - no → leaf 20: 4.145
      - no → Duration_x_Heart_Rate ≤1601.000
        - yes → Age_plus_Heart_Rate ≤133.500
          - yes → leaf 7: 4.145
          - no → leaf 22: 4.148
        - no → Age_plus_Heart_Rate ≤137.500
          - yes → leaf 15: 4.149
          - no → leaf 23: 4.152
    - no → Duration_x_Heart_Rate ≤2551.000
      - yes → Duration_x_Heart_Rate ≤2182.000
        - yes → Age_plus_Heart_Rate ≤140.500
          - yes → leaf 3: 4.153
          - no → leaf 27: 4.156
        - no → Age_plus_Heart_Rate ≤140.500
          - yes → leaf 14: 4.156
          - no → leaf 24: 4.159
      - no → Duration_plus_Heart_Rate ≤135.500
        - yes → Age_x_Weight ≤3720.500
          - yes → leaf 6: 4.161
          - no → leaf 21: 4.164
        - no → row_mean ≤82.925
          - yes → leaf 13: 4.164
          - no → leaf 28: 4.168

---

**6. Challenges Faced**

**Initial Issues**:

- Overfitting due to extensive feature engineering, mitigated by cross-validation.

- Slow training times for CatBoost, addressed by optimizing early stopping.

**Tradeoffs and Limitations**:

- **Computational Cost**: Extensive feature engineering and ensemble training required significant resources.

- **Generalization**: The model may not generalize well to datasets with different feature distributions.

- **Lack of Real-time UI**: The notebook-based interface limits practical deployment.