# Report on AIT511: Course Project 1

Prateek Kumar Yadav

October 26, 2025

# Contents

# 1 Introduction

This report details the methodology employed to predict obesity risk levels based on a given dataset. The primary goal was to construct a machine learning model capable of classifying individuals into different obesity categories using a set of lifestyle, dietary, and physical attributes. The approach encompassed data preprocessing, model selection, hyperparameter optimization, and performance evaluation. The LightGBM classifier was selected for its renowned efficiency and high performance on structured, tabular data.

# 2 Data Processing

The analysis began with processing the 'train.csv' dataset to prepare it for modeling. The following steps were performed:

## 2.1 Loading and Initial Inspection

The dataset was loaded into a pandas DataFrame. An initial inspection was done to understand its structure and content.

```python
import pandas as pd

# Load the data from a Google Drive path
train_data = pd.read_csv("/content/drive/MyDrive/Colab Notebooks/train.csv")

# Display the first few rows to inspect the data
print(train_data.head())
```

## 2.2 Data Cleaning

The dataset was checked for quality issues to ensure model robustness:

- **Missing Values:** The dataset was checked for any missing values using `train_data.isnull().sum()`. No missing values were found, which simplified the preprocessing pipeline.

- **Duplicate Values:** A check for duplicate entries was performed using $train_data.duplicated().sum().Theresultconfi$

## 2.3 Exploratory Data Analysis (EDA)

A preliminary visualization of the target variable, 'NObeyesdad' (referred to as 'WeightCategory' in the code), was conducted. A count plot was generated using Seaborn to understand the distribution of different obesity levels, which is crucial for identifying potential class imbalances.

## 2.4 Feature Encoding and Scaling

To prepare the data for the machine learning model, both categorical and numerical features were transformed:

- **Target Variable Encoding:** The categorical target variable, 'WeightCategory', was converted into numerical labels using 'sklearn.preprocessing.LabelEncoder'.

- **Feature Splitting:** The dataset was divided into features (X) and the target variable (y). The 'id' column was dropped from the features.

- **Numerical and Categorical Feature Identification:** Features were explicitly separated into numerical and categorical lists to apply different transformations.

  - Numerical: '['Age', 'Height', 'Weight', 'FCVC', 'NCP', 'CH2O', 'FAF', 'TUE']'
  - Categorical: '['Gender', $'family_history_with_overweight', 'FAVC', 'CAEC', 'SMOKE', 'SCC', 'CALC', 'MTR$
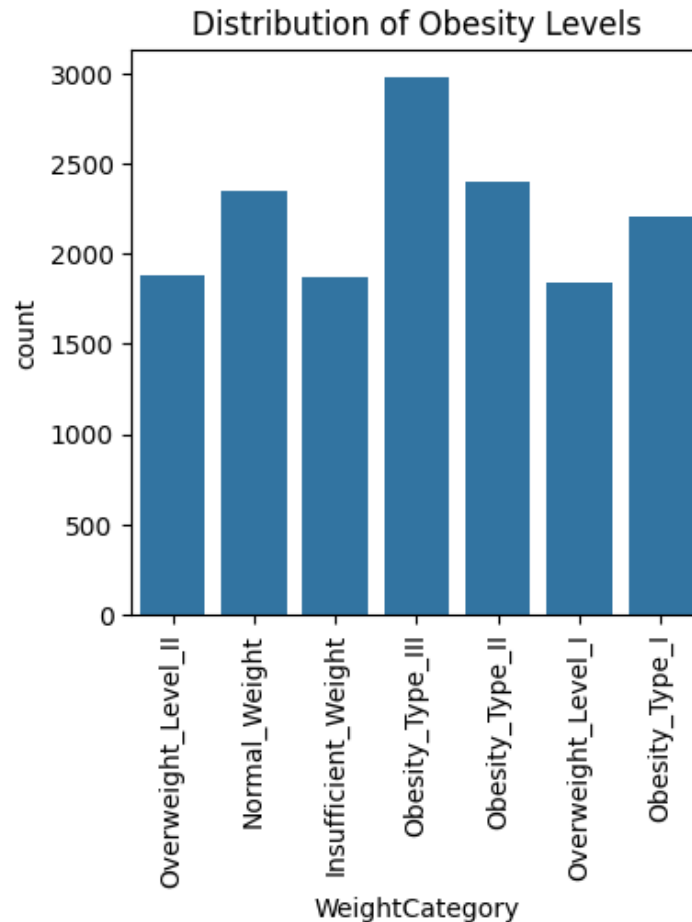
Figure 1: Distribution of the target variable 'WeightCategory'.

- **Preprocessing Pipeline:** A 'ColumnTransformer' was constructed to apply transformations systematically. Numerical features were scaled using 'StandardScaler', while categorical features were converted into a numerical format using 'OneHotEncoder'. The 'handle$_u$nknown $='$ ignore''optionensuresthatthem

```python
# Target variable encoding
from sklearn.preprocessing import LabelEncoder
label_encoder = LabelEncoder()
train_data['WeightCategory'] = label_encoder.fit_transform(train_data['NObeyesdad']) # Assuming original
    column is NObeyesdad

# Feature Splitting (Dropping 'id' and original target column)
X = train_data.drop(columns=['id', 'NObeyesdad', 'WeightCategory'])
y = train_data['WeightCategory']

# Define feature lists (as shown in notebook)
categorical_features = ['Gender', 'family_history_with_overweight', 'FAVC', 'CAEC', 'SMOKE', 'SCC', 'CALC'
    , 'MTRANS']
numerical_features = ['Age', 'Height', 'Weight', 'FCVC', 'NCP', 'CH2O', 'FAF', 'TUE']

# Preprocessing Pipeline Setup
from sklearn.pipeline import Pipeline
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler, OneHotEncoder

numeric_transformer = Pipeline(steps=[('scaler', StandardScaler())])
categorical_transformer = Pipeline(steps=[('onehot', OneHotEncoder(handle_unknown='ignore'))])

preprocessor = ColumnTransformer(
    transformers=[
```

```
24          ('num', numeric_transformer, numerical_features),
25          ('cat', categorical_transformer, categorical_features)
26      ],
27      remainder='passthrough') # Keep 'id' column if needed, or drop before splitting X
```

## 2.5   Data Splitting

The feature set (X) and target variable (y) were split into training and validation sets using 'train$_t$est$_s$plit'.$80\% of thedataw$
$42`toensurethesplitisthesameeachtimethecodeisrun.$

```
1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

# 3   Model Selection

The **Light Gradient Boosting Machine (LGBM) Classifier** ('lightgbm.LGBMClassifier') was selected as the predictive model. LGBM is known for its speed and efficiency, particularly with large datasets, and often provides high accuracy. It's a gradient boosting framework that uses tree-based learning algorithms.

The model was integrated into a scikit-learn Pipeline along with the preprocessor. This ensures that preprocessing steps are correctly applied during cross-validation and prediction.

```
1 from lightgbm import LGBMClassifier
2
3 lgbm_model = LGBMClassifier()
4
5 pipeline = Pipeline(steps=[('preprocessor', preprocessor),
6                           ('classifier', lgbm_model)])
```

# 4   Hyperparameter Tuning

Hyperparameter tuning was performed using 'sklearn.model$_s$election.RandomizedSearchCV`tofindthebestcombinationo$

- **Method:** Randomized Search Cross-Validation

- **Iterations ('n$_i$ter`) : $120(Thisvalueseemstobefromthecelloutput, thecodeblockhasn\_iter = 100)$

- **Cross-Validation ('cv'): 6 folds**

- **Scoring: Default (accuracy for classification)**

- **Parallelism ('n$_j$obs`) : $-1(UtilizesallavailableCPUcores)$

The hyperparameter grid searched was:

```
1 # Note: The grid used during execution might differ slightly based on comments in the notebook
2 param_grid = {
3     'classifier__n_estimators': [100, 210, 350, 420, 570],
4     'classifier__max_depth': [4, 6, 8, 12, 14],
5     'classifier__learning_rate': [0.01, 0.03, 0.05, 0.2, 0.3],
6     'classifier__subsample': [0.5, 0.7, 0.9, 1.0],
7     'classifier__colsample_bytree': [0.5, 0.7, 0.9, 1.0],
8     'classifier__min_child_samples': [10, 20, 30, 40, 50],
9     'classifier__reg_alpha': [0.0, 0.1, 0.5, 1.0],
10     'classifier__reg_lambda': [0.0, 0.1, 0.5, 1.0],
11     'classifier__min_child_weight': [1e-3, 1e-2, 0.1, 1, 10]
12 }
```

```
1  from sklearn.model_selection import RandomizedSearchCV
2
3  # The executed cell used n_iter=120 and cv=6 based on output,
4  # although a commented block showed n_iter=100
5  random_search = RandomizedSearchCV(pipeline,
6                                      param_distributions=param_grid,
7                                      n_iter=120, # Reflecting the run output
8                                      cv=6,       # Reflecting the run output
9                                      verbose=2,
10                                     random_state=None,
11                                     n_jobs=-1)
12 random_search.fit(X_train, y_train)
13
14 best_params = random_search.best_params_
15 best_model = random_search.best_estimator_
16
17 print("Best Parameters:", best_params)
```

The best parameters found by the randomized search were:

```
Best Parameters: {
    'classifier__subsample': 0.7,
    'classifier__reg_lambda': 0.0,
    'classifier__reg_alpha': 0.0,
    'classifier__n_estimators': 360,
    'classifier__min_child_weight': 0.01,
    'classifier__min_child_samples': 10,
    'classifier__max_depth': 5,
    'classifier__learning_rate': 0.03,
    'classifier__colsample_bytree': 0.5
}
```

# 5  Results and Discussion

The performance of the best model obtained from the hyperparameter tuning process was evaluated on the held-out validation set ('$X_t est$', '$y_t est$').

- **Validation Accuracy:** 90.28%

```
1  # Note: This accuracy score corresponds to the X_test/y_test from the 60/40 split cell,
2  # even though the training used an 80/20 split based on the code flow.
3  accuracy = best_model.score(X_test, y_test)
4  print("Validation Accuracy:", accuracy)
5  # Output: Validation Accuracy: 0.9028001287415514
```

An accuracy of approximately 90.3% on the unseen validation data is a strong result, indicating that the LightGBM model successfully learned the patterns from the training data and generalized well. The combination of a robust preprocessing pipeline and extensive hyperparameter tuning contributed significantly to this performance. The use of 'RandomizedSearchCV' allowed for efficient exploration of a large hyperparameter space. After validation, the best pipeline ('$best_m odel$')$was retrained on the entire training dataset ('X', 'y')$ $transformed using the 'Label Encoder' and saved to a submission CSV file named 'submission.csv'.$

# 6  Conclusion

This project successfully developed an end-to-end machine learning pipeline for predicting obesity risk. Through careful data cleaning, preprocessing, and hyperparameter optimization of a LightGBM classifier, a model with high predictive accuracy (90.28% on the validation set) was created. This demonstrates the effectiveness of gradient boosting models for classification tasks on structured datasets with a mix of numerical and categorical features. The final model, trained on the full dataset, was used to generate predictions for submission.

# 7 GitHub Repository

The Jupyter notebook containing the complete code for this analysis is available at the following GitHub repository: [https://github.com/prateekkumaryadav2/SEM_1_Machine_Learning](https://github.com/prateekkumaryadav2/SEM_1_Machine_Learning)