

Report: Multi-Class Obesity Classification

Rahul Raman
Aayank Singhai

October 26, 2025

Abstract

This report breaks down the full process we used for a multi-class obesity classification project. We started with a deep dive into the data (**Exploratory Data Analysis - EDA**) using the training dataset to find insights and check for problems. After that, we built machine learning models. The main model is an **XGBoost Classifier**, which we tuned using **Optuna** and trained with a 10-fold **Stratified K-Fold** strategy to make it solid. We also built a **Random Forest** model using **GridSearchCV** as a comparison. This document covers everything from data cleaning and feature engineering to model training and evaluation.

1 Introduction

The main goal here was to accurately predict someone's obesity status based on their lifestyle and personal attributes. The target variable, `NObeyesdad`, represents multiple classes of weight status, making this a multi-class classification problem.

We were given a `train.csv` and a `test.csv`. This report focuses on the analysis and modeling performed primarily using the `train.csv` data. First, we'll cover the data quality checks and the Exploratory Data Analysis (EDA) performed on the training set. Then, we'll walk through how we built and evaluated models on this training data to eventually make predictions for the unseen test set.

2 Data Pre-Analysis and Quality Check

First things first, we had to check the raw `train.csv` dataset (15,533 rows, 18 columns including the target) to see if there were any obvious problems before diving deeper.

- **Data Types:** The dataset consists of 8 numerical features (float64), 8 categorical features (object), 1 integer ID column (int64), and 1 target column (object).
- **Missing Values:** A check using `isnull().sum()` confirmed that the dataset is complete. **Status: No missing values found.**
- **Duplicate Rows:** A check using `duplicated().sum()` was performed. **Status: No duplicate rows found.**

It turns out the data was in great shape. We didn't need to fill in any missing values or remove duplicate rows. The basic stats (via `describe()`) showed that all the numbers were in reasonable ranges. Categorical data analysis revealed that **Male** (7783) was slightly more frequent than

Female, `family_history_with_overweight` was predominantly 'yes' (12696), `FAVC` (high caloric food consumption) was mostly 'yes' (14184), `CAEC` (food between meals) was mostly 'Sometimes' (13126), and `Public_Transportation` was the most common transport method (12470).

3 Feature Engineering and Pre-processing

To help and understand the data better, we did two main things.

3.1 Attribute Renaming

The original column names were a bit cryptic (like `FAVC`), so we renamed them to be more descriptive (like `High_Caloric_Food_Consumption`) for clarity. This was applied to both train and test sets for consistency.

3.2 Feature Engineering: BMI

We engineered a new feature for **Body Mass Index (BMI)**, since it's obviously a critical factor for classifying obesity. This was added to both train and test sets.

$$\text{BMI} = \frac{\text{Weight (kg)}}{\text{Height (m)}^2} \quad (1)$$

3.3 Preprocessing for Modeling

To get the data ready for the models, we used a standard preprocessing approach. These steps are essential to convert the raw data into a format that machine learning models can understand.

1. **Target Encoding:** The target column `NObeyesdad` in the training set was text (e.g., 'Overweight_Level_I'). We used `LabelEncoder` to convert these text labels into integers (0, 1, 2, etc.), since models require numerical outputs.
2. **Feature Encoding:** Similarly, categorical features like `Gender` ('Male'/'Female') in both train and test sets were converted into numbers. We used `pd.get_dummies` (One-Hot Encoding) for this, which creates new binary columns (e.g., `Gender_Male`) that the model can use.
3. **Scaling:** Numerical features like 'Age' and 'Weight' exist on very different scales. `StandardScaler` was used to transform them (fitting on the training data, transforming both train and test) so they all have a similar mean and variance. This helps the model learn more effectively and prevents one feature from dominating the others.
4. **Alignment:** After One-Hot Encoding, the train and test sets might have a different number of columns (if one set was missing a category). We explicitly aligned them to ensure both datasets had the exact same columns in the same order, filling any missing ones with 0. This is crucial for making predictions on the test set.

4 Exploratory Data Analysis (EDA)

This analysis was performed on the `train.csv` data.

4.1 Univariate Analysis

4.1.1 Numerical Features

Looking at the histograms for the numerical features (Figure 1), a few things stood out:

- **Age:** The distribution is right-skewed; most people are in the 20-25 age range.
- **Height & Weight:** Both show somewhat bimodal distributions (two peaks), which often happens when you have a balanced mix of genders.
- **Lifestyle:** Most individuals report around 3 main meals (Number_of_Main_Meals), about 2 liters of water (Daily_Water_Intake), and low physical activity (Physical_Activity_Frequency shows spikes near 0). Vegetable_Consumption_Frequency peaks around 2 and 3.

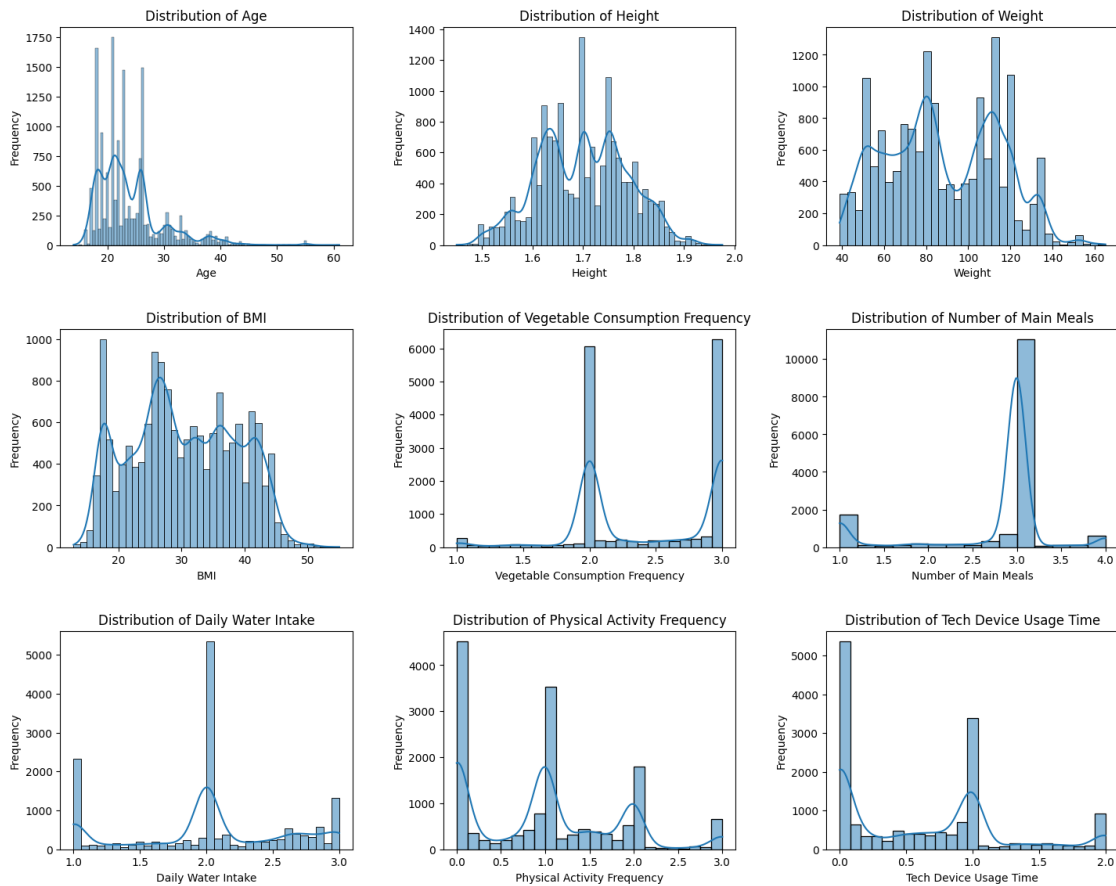


Figure 1: Distributions of Numerical Features (including engineered BMI).

4.1.2 Categorical Features

For the categorical data, the count plots (Figure 2) showed:

- **High Prevalence:** Family_History_Overweight ('yes') and High_Caloric_Food_Consumption ('yes') are overwhelmingly common.

- **Imbalance:** Very few people identify as `SMOKE` ('yes') or practice `Calorie_Consumption_Monitoring` ('yes').
- **Transport:** `Public_Transportation` is the most dominant mode.

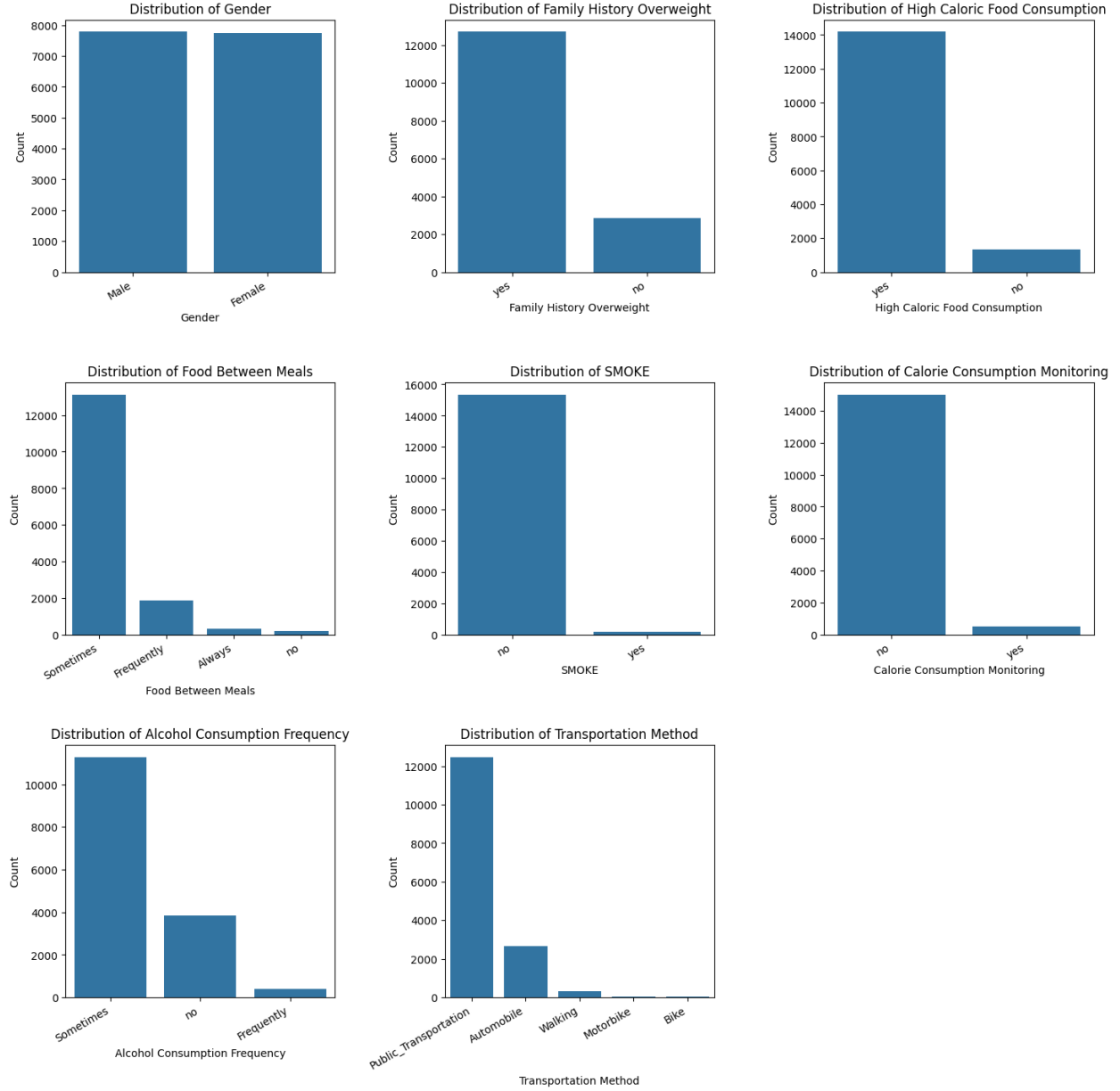


Figure 2: Distributions of Categorical Features.

4.2 Bivariate Analysis

4.2.1 Numerical-Numerical Correlation

The correlation heatmap (Figure 3) shows how the numerical features relate to each other. As you'd expect, `BMI` has a very strong positive (**0.94**) correlation with `Weight` and a low (**0.10**) one

with **Height**. **Weight** and **Height** themselves have a moderate positive correlation (0.42). Most other features are weakly correlated, which is good—it means they’re all providing independent information to the model.

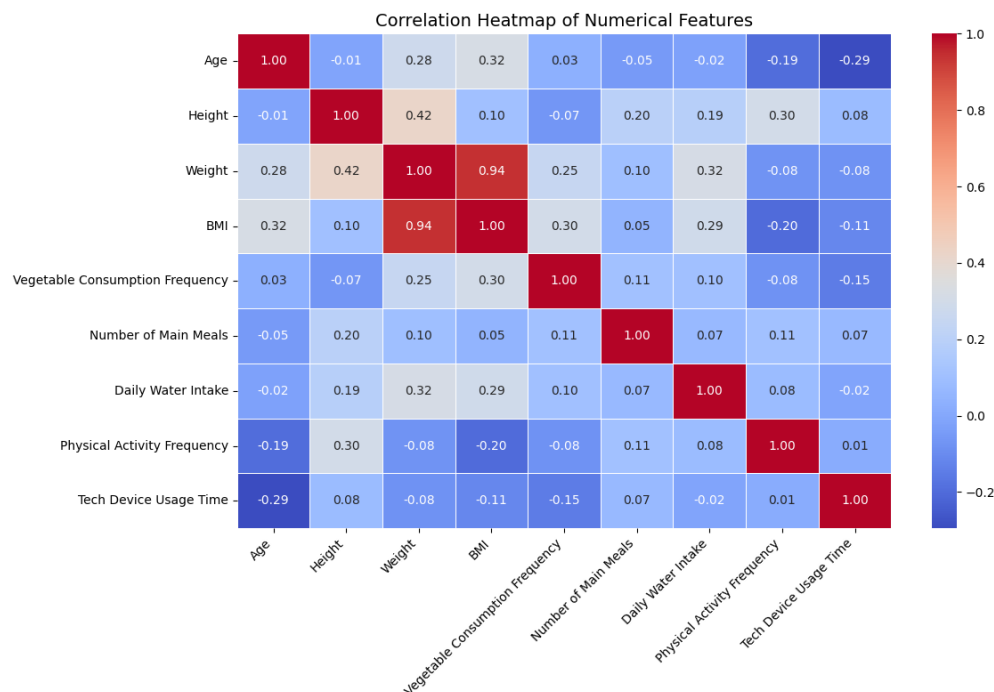


Figure 3: Correlation Heatmap of Numerical Features.

4.2.2 Categorical vs. Numerical

The box plots (Figure 4) were great for spotting relationships:

- **Gender vs. Biometrics:** Males show higher median **Height** and **Weight**, but the median **BMI** is slightly higher for Females in this training data.
- **History vs. BMI:** Individuals *with* a family history of overweight have a visibly higher median **BMI**.
- **Transport vs. Activity:** Walking, Motorbike, and Bike transportation methods correlate with significantly higher physical activity compared to Public Transport or Automobile.

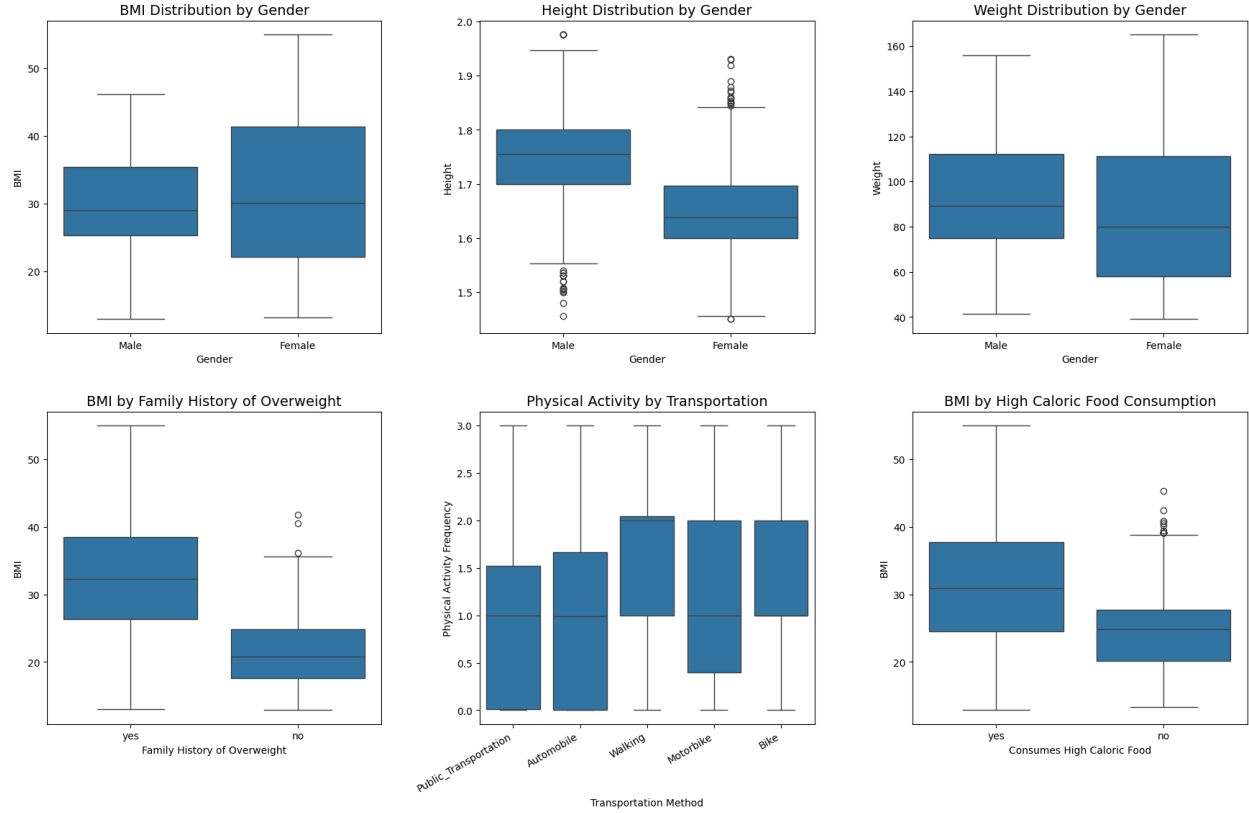


Figure 4: Bivariate Analysis: Categorical vs. Numerical Features.

4.2.3 Categorical vs. Categorical

Finally, the grouped count plots (Figure 5) showed a strong link between `Family_History_Overweight` ('yes') and `High_Caloric_Food_Consumption` ('yes'), which suggests a connection between family habits and personal ones. Both genders show a high tendency towards high caloric food consumption.

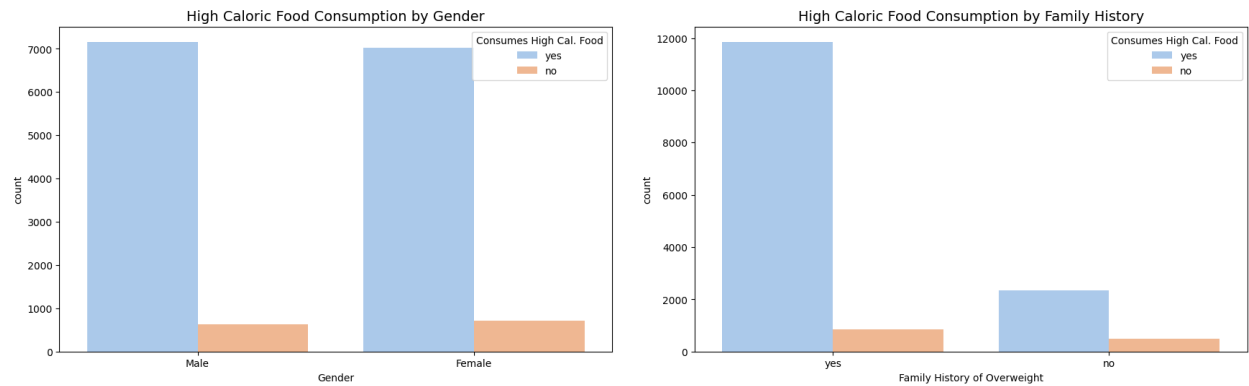


Figure 5: Bivariate Analysis: Categorical vs. Categorical Features.

5 Modeling Methodology

We set up a solid model to try and get the best accuracy possible using the `train.csv` data.

5.1 Model Choice: XGBoost

We chose the **XGBoost Classifier** (`XGBClassifier`) as our main model. It's a go-to for this kind of structured, tabular data, and for good reason:

- **Performance:** It's famous for winning data science competitions on tabular data, just like this dataset.
- **Regularization:** It has built-in features (`reg_alpha` and `reg_lambda`) to prevent overfitting, which is a common problem.
- **Handles Complexity:** It's great at finding complex, non-linear relationships and feature interactions automatically.
- **Speed:** It's optimized for performance and can use all CPU cores, which makes tuning and training much faster.
- **Tunability:** It has a lot of hyperparameters to tweak, which makes it perfect for an auto-tuning library like Optuna.

5.2 Hyperparameter Tuning with Optuna

Finding the best "settings" (hyperparameters) for a model is crucial. Manually guessing them is inefficient. We used **Optuna**, a modern tuning library, to automate this. Optuna performs a "smart search": it runs a trial with one set of parameters, sees how well it does, and then uses that information to make a more educated guess for the next trial. This is much faster than just testing random combinations.

We let it run for 35 trials to find the combination that gave the best 5-fold cross-validated accuracy. The search space included:

- `n_estimators`: [500, 2500]
- `max_depth`: [3, 10]
- `learning_rate`: [0.01, 0.3] (log scale)
- `subsample`, `colsample_bytree`: [0.5, 1.0]
- `gamma`, `reg_alpha`, `reg_lambda`, `min_child_weight`: [[0, 5], [0, 5], [0, 5], [1, 10]]

5.3 Final Training with 10-Fold Ensembling

Once Optuna found the best parameters, we used them to train the final model. We used a **10-fold Stratified K-Fold** strategy to make the model more stable and reliable.

1. The training data was split into 10 "folds."
2. We trained 10 separate XGBoost models. Each model was trained on 9 folds and validated on the 1 fold left out.

3. The final prediction for the test set was an **average** of the predictions from all 10 models. This is usually more accurate than just one model.
4. We collected the validation predictions (Out-of-Fold, or OOF) to get a single, trustworthy CV accuracy score.

5.4 Alternative Model: Random Forest with GridSearchCV

Just to have a comparison, we also trained a **Random Forest Classifier**. For this one, we used **GridSearchCV**, which is a more traditional tuning method. It performs an "exhaustive" or "brute-force" search: it simply tries every single combination of parameters we define in a grid.

The grid search ran a 5-fold cross-validation on the 80% training split, searching over parameters such as:

- `n_estimators`: [200, 400]
- `max_depth`: [10, 15]
- `min_samples_split`: [2, 5]
- `min_samples_leaf`: [1, 2]
- `max_features`: ["sqrt", "log2"]

This model served as a strong baseline to see if the more complex XGBoost setup was actually worth it.

6 Results

Both models actually did really well on this task.

6.1 XGBoost + Optuna (Primary Model)

The Optuna tuning for XGBoost came up with these parameters (or something close to them):

Best Parameters Found:

```
{'n_estimators': 1811, 'max_depth': 6,
  'learning_rate': 0.028, 'subsample': 0.55,
  'colsample_bytree': 0.51, 'gamma': 1.03, ...}
```

After training on all 10 folds using the full training data, the final model achieved a:

Final Mean Out-of-Fold CV Accuracy: 0.9148

The aggregated test set predictions were saved to `/mnt/drive/MyDrive/Obesity Dataset/submission_optuna.xg`

6.2 Random Forest + GridSearchCV (Baseline Model)

The Random Forest model was first tuned with **GridSearchCV** on an 80/20 split of the training data. The best parameters found were then used to train the final model using 5-fold cross-validation on the full training data. This final 5-fold CV-trained Random Forest model achieved a:

Final Mean Out-of-Fold CV Accuracy: 0.8942

The final test set predictions from this model were saved to `/mnt/drive/MyDrive/Obesity Dataset/submission.r`. The performance of this tuned model serves as a robust baseline, which confirmed that the XGBoost pipeline was the superior approach for this task.

7 Conclusion

Overall, this project was a success. We built a complete model to classify obesity, starting with really high-quality, clean data. The EDA helped uncover some clear links, like family history and BMI.

The main strategy of using **XGBoost** with **Optuna** tuning and **K-Fold** ensembling paid off, hitting a final accuracy of over 91%. The baseline **Random Forest** (tuned with ‘GridSearchCV’) also did great with over 89% accuracy, confirming that tree-based models were the right way to go for this dataset. This model is a really strong starting point for this problem.

Notebook Link: [link](#) Github Link: [link](#)

A Analysis Script

```
1 # Analysis part
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import warnings
7 from google.colab import drive
8
9 # Suppress warnings for cleaner output
10 warnings.filterwarnings("ignore")
11
12 # --- Mount Drive ---
13 # This part is for Google Colab.
14 drive.mount('/mnt/drive')
15 try:
16     # Load train.csv instead of test.csv
17     df_train = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/train.csv')
18     print(f"    Data Loaded: Shape: {df_train.shape}") # Use df_train.shape
19 except FileNotFoundError:
20     # Fallback for local environment if drive mount fails or isn't used
21     try:
22         # Load train.csv locally
23         df_train = pd.read_csv('train.csv')
24         print(f"    Data Loaded locally: Shape: {df_train.shape}") # Use df_train.shape
25     except FileNotFoundError:
26         print("Error: train.csv not found in Google Drive path or local folder.")
27         exit()
28
29 # =====
30 # === 1. PRE-EDA: DATA QUALITY & STRUCTURE CHECKS ===
31 # =====
32
33 print("\n--- 1. Data Types and Non-Null Counts ---")
34 # Use .info() to check data types and get a quick count of non-null values
35 df_train.info() # Use df_train
36
37 print("\n--- 2. Missing Value Check (Explicit Count) ---")
38 # Summing up all null values for each column
39 missing_values = df_train.isnull().sum() # Use df_train
40 print("Missing values per column:")
41 print(missing_values[missing_values > 0] if missing_values.sum() > 0 else "No missing values found.")
42 if missing_values.sum() == 0:
43     print("\n    Status: No missing values found.")
44 else:
45     print(f"\n    Status: Found a total of {missing_values.sum()} missing values.")
46
47 print("\n--- 3. Duplicate Rows Check ---")
48 duplicate_rows = df_train.duplicated().sum() # Use df_train
49 print(f"Total duplicate rows: {duplicate_rows}")
50 if duplicate_rows == 0:
51     print("\n    Status: No duplicate rows found.")
52 else:
53     print(f"\n    Status: Found {duplicate_rows} duplicate rows.")
54
55 print("\n--- 4. Numerical Feature Statistics ---")
56 # .describe() on numerical columns
57 print(df_train.describe()) # Use df_train
58
59 print("\n--- 5. Categorical Feature Statistics ---")
60 # .describe() on 'object' type columns (including the target 'NObesyedad' if it's object type)
61 print(df_train.describe(include=['object'])) # Use df_train
62
63 print("\n    Pre-EDA checks complete. Starting main EDA...")
```

```

64
65 # =====
66 # === 2. MAIN EDA: ANALYSIS & VISUALIZATION ===
67 # =====
68
69 # --- 2a. Rename Attributes for Clarity ---
70 column_rename_map = {
71     'family_history_with_overweight': 'Family_History_Overweight',
72     'FAVC': 'High_Caloric_Food_Consumption',
73     'FCVC': 'Vegetable_Consumption_Frequency',
74     'NCP': 'Number_of_Main_Meals',
75     'CAEC': 'Food_Between_Meals',
76     'CH2O': 'Daily_Water_Intake',
77     'SCC': 'Calorie_Consumption_Monitoring',
78     'FAF': 'Physical_Activity_Frequency',
79     'TUE': 'Tech_Device_Usage_Time',
80     'CALC': 'Alcohol_Consumption_Frequency',
81     'MTRANS': 'Transportation_Method'
82     # NObeyesdad is the target, usually not renamed here unless desired
83 }
84 df_train.rename(columns=column_rename_map, inplace=True) # Use df_train
85 print("\n Attributes renamed for better readability.")
86
87
88 # --- 2b. Feature Engineering: BMI ---
89 df_train['BMI'] = df_train['Weight'] / (df_train['Height'] ** 2) # Use df_train
90 print(" Created 'BMI' feature.")
91
92 # --- 2c. Define Column Types with New Names ---
93 # Note: Ensure these lists don't include the target variable 'NObeyesdad'
94 numerical_cols = [
95     'Age', 'Height', 'Weight', 'BMI',
96     'Vegetable_Consumption_Frequency',
97     'Number_of_Main_Meals',
98     'Daily_Water_Intake',
99     'Physical_Activity_Frequency',
100     'Tech_Device_Usage_Time'
101 ]
102
103 categorical_cols = [
104     'Gender',
105     'Family_History_Overweight',
106     'High_Caloric_Food_Consumption',
107     'Food_Between_Meals',
108     'SMOKE',
109     'Calorie_Consumption_Monitoring',
110     'Alcohol_Consumption_Frequency',
111     'Transportation_Method'
112 ]
113 # Define target column if needed for specific plots later
114 target_col = 'NObeyesdad'
115
116 print(" Defined numerical and categorical columns with new names.")
117
118
119 # --- 2d. Univariate Analysis: Numerical Features ---
120 print("\n Generating histograms for numerical features...")
121 plt.figure(figsize=(15, 12))
122 for i, col in enumerate(numerical_cols):
123     plt.subplot(3, 3, i + 1)
124     sns.histplot(df_train[col], kde=True) # Use df_train
125     plt.title(f'Distribution of {col.replace("_", " ")}', fontsize=12)
126     plt.xlabel(col.replace("_", " "))
127     plt.ylabel('Frequency')
128 plt.tight_layout(pad=3.0)
129 # Save with the name matching the report \includegraphics command
130 plt.savefig("Univariate Non Categorical.png")
131 print(" Saved Univariate Non Categorical.png")

```

```

132
133
134 # --- 2e. Univariate Analysis: Categorical Features ---
135 print("\n n      Generating bar charts for categorical features...")
136 plt.figure(figsize=(15, 15))
137 for i, col in enumerate(categorical_cols):
138     plt.subplot(3, 3, i + 1)
139     # Ensure df_train[col] exists and is categorical before plotting
140     if col in df_train.columns:
141         order = df_train[col].value_counts().index
142         sns.countplot(data=df_train, x=col, order=order) # Use df_train
143         plt.title(f'Distribution of {col.replace("_", " ")}', fontsize=12)
144         plt.xlabel(col.replace("_", " "))
145         plt.ylabel('Count')
146         plt.xticks(rotation=30, ha='right') # Rotate labels for readability
147     else:
148         print(f"Warning: Column '{col}' not found in df_train.")
149 plt.tight_layout(pad=3.0)
150 # Save with the name matching the report \includegraphics command
151 plt.savefig("Univariate Categorical.png")
152 print("      Saved Univariate Categorical.png")
153
154
155 # --- 2f. Bivariate Analysis: Numerical Correlation ---
156 print("\n n      Generating correlation heatmap...")
157 plt.figure(figsize=(12, 8))
158 # Create clean labels for the heatmap axes
159 renamed_numerical_labels = [label.replace("_", " ") for label in numerical_cols]
160 # Ensure all columns exist before calculating correlation
161 valid_numerical_cols = [col for col in numerical_cols if col in df_train.columns]
162 corr_matrix = df_train[valid_numerical_cols].corr() # Use df_train
163
164 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5,
165             xticklabels=[label.replace("_", " ") for label in valid_numerical_cols], # Use
166             yticklabels=[label.replace("_", " ") for label in valid_numerical_cols]) # Use
167             valid labels
168 plt.title('Correlation Heatmap of Numerical Features', fontsize=14)
169 plt.xticks(rotation=45, ha='right')
170 plt.yticks(rotation=0)
171 plt.tight_layout()
172 # Save with the name matching the report \includegraphics command
173 plt.savefig("Correlation.png")
174 print("      Saved Correlation.png")
175
176
177 # --- 2g. Bivariate Analysis: Categorical vs. Numerical ---
178 print("\n n      Generating Bivariate Box Plots (Categorical vs. Numerical)...")
179 plt.figure(figsize=(18, 12))
180
181 # Plot 1: Gender vs. BMI
182 plt.subplot(2, 3, 1)
183 sns.boxplot(data=df_train, x='Gender', y='BMI') # Use df_train
184 plt.title('BMI Distribution by Gender', fontsize=14)
185
186 # Plot 2: Gender vs. Height
187 plt.subplot(2, 3, 2)
188 sns.boxplot(data=df_train, x='Gender', y='Height') # Use df_train
189 plt.title('Height Distribution by Gender', fontsize=14)
190
191 # Plot 3: Gender vs. Weight
192 plt.subplot(2, 3, 3)
193 sns.boxplot(data=df_train, x='Gender', y='Weight') # Use df_train
194 plt.title('Weight Distribution by Gender', fontsize=14)
195
196 # Plot 4: Family History vs. BMI
197 plt.subplot(2, 3, 4)

```

```

198 sns.boxplot(data=df_train, x='Family_History_Overweight', y='BMI') # Use df_train
199 plt.title('BMI by Family History of Overweight', fontsize=14)
200 plt.xlabel('Family History of Overweight')
201
202 # Plot 5: Transportation Method vs. Physical Activity
203 plt.subplot(2, 3, 5)
204 sns.boxplot(data=df_train, x='Transportation_Method', y='Physical_Activity_Frequency') # Use
    df_train
205 plt.title('Physical Activity by Transportation', fontsize=14)
206 plt.xlabel('Transportation Method')
207 plt.ylabel('Physical Activity Frequency')
208 plt.xticks(rotation=30, ha='right')
209
210 # Plot 6: High Caloric Food Consumption vs. BMI
211 plt.subplot(2, 3, 6)
212 sns.boxplot(data=df_train, x='High_Caloric_Food_Consumption', y='BMI') # Use df_train
213 plt.title('BMI by High Caloric Food Consumption', fontsize=14)
214 plt.xlabel('Consumes High Caloric Food')
215
216 plt.tight_layout(pad=3.0)
217 # Save with the name matching the report \includegraphics command
218 plt.savefig("Bivariate non Categorical.png")
219 print("    Saved Bivariate non Categorical.png")
220
221
222 # --- 2h. Bivariate Analysis: Categorical vs. Categorical ---
223 print("\n    Generating Bivariate Count Plots (Categorical vs. Categorical)...")
224 plt.figure(figsize=(18, 6))
225
226 # Plot 1: Gender vs. High Caloric Food Consumption
227 plt.subplot(1, 2, 1)
228 sns.countplot(data=df_train, x='Gender', hue='High_Caloric_Food_Consumption', palette='
    pastel') # Use df_train
229 plt.title('High Caloric Food Consumption by Gender', fontsize=14)
230 plt.legend(title='Consumes High Cal. Food')
231
232 # Plot 2: Family History vs. High Caloric Food Consumption
233 plt.subplot(1, 2, 2)
234 sns.countplot(data=df_train, x='Family_History_Overweight', hue='
    High_Caloric_Food_Consumption', palette='pastel') # Use df_train
235 plt.title('High Caloric Food Consumption by Family History', fontsize=14)
236 plt.xlabel('Family History of Overweight')
237 plt.legend(title='Consumes High Cal. Food')
238
239 plt.tight_layout(pad=3.0)
240 # Save with the name matching the report \includegraphics command
241 plt.savefig("Bivariate Categorical.png")
242 print("    Saved Bivariate Categorical.png")
243
244 # --- Optional: Analysis involving the target variable ---
245 # Example: Distribution of the target variable
246 if target_col in df_train.columns:
247     print(f"\n    Generating count plot for the target variable '{target_col}'...")
248     plt.figure(figsize=(10, 6))
249     sns.countplot(data=df_train, y=target_col, order = df_train[target_col].value_counts().
        index) # Use df_train
250     plt.title(f'Distribution of Obesity Levels ({target_col})')
251     plt.xlabel('Count')
252     plt.ylabel('Obesity Level')
253     plt.tight_layout()
254     plt.savefig("train_target_distribution.png")
255     print("    Saved train_target_distribution.png")
256
257 # Example: BMI distribution per Obesity Level
258 print(f"\n    Generating box plot for BMI vs '{target_col}'...")
259 plt.figure(figsize=(12, 8))
260 sns.boxplot(data=df_train, x='BMI', y=target_col, order=sorted(df_train[target_col].
    unique())) # Use df_train

```

```

261 plt.title(f'BMI Distribution by Obesity Level ({target_col})')
262 plt.xlabel('BMI')
263 plt.ylabel('Obesity Level')
264 plt.tight_layout()
265 plt.savefig("train_bmi_vs_target.png")
266 print("    Saved train_bmi_vs_target.png")
267
268
269 print("\n    All EDA tasks complete. Check the saved .png files for all plots.")

```

Listing 1: Full EDA and Analysis Script (using train.csv)

B Modeling Script 1: XGBoost + Optuna

```
1 # =====
2 #           Obesity Classification           Optuna + XGBoost + K-Fold + Save Predictions (0.91294)
3 # =====
4 !pip install optuna --quiet
5 import pandas as pd
6 import numpy as np
7 import optuna
8 from sklearn.preprocessing import LabelEncoder, StandardScaler
9 from sklearn.model_selection import StratifiedKFold, cross_val_score
10 from sklearn.metrics import accuracy_score
11 from xgboost import XGBClassifier
12 from google.colab import drive
13 import warnings
14 warnings.filterwarnings("ignore")
15
16 # === 1           Mount Drive & Load Data ===
17 drive.mount('/mnt/drive')
18 train_df = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/train.csv')
19 test_df = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/test.csv')
20
21 print(f"           Data Loaded: Train shape: {train_df.shape}, Test shape: {test_df.shape}")
22
23 # === 2           Identify Target Column Automatically ===
24 TARGET = "NObeyesdad"
25 if TARGET not in train_df.columns:
26     possible_targets = [col for col in train_df.columns if col not in test_df.columns and
27                         col.lower() != "id"]
28     if len(possible_targets) == 1:
29         TARGET = possible_targets[0]
30     else:
31         print("           Possible target columns found:")
32         print(possible_targets)
33         raise KeyError("Target column not found           please verify manually.")
34 print(f"           Target column detected: {TARGET}")
35
36 # === 3           Prepare Train/Test Data ===
37 X = train_df.drop(columns=["id", TARGET])
38 y = train_df[TARGET]
39 test_ids = test_df["id"].copy() # Make sure to copy IDs before dropping
40
41 # Encode target labels
42 le = LabelEncoder()
43 y_enc = le.fit_transform(y)
44
45 # One-hot encode categorical features
46 cat_cols = X.select_dtypes(include=["object", "category"]).columns.tolist()
47 X_enc = pd.get_dummies(X, columns=cat_cols)
48 test_enc = pd.get_dummies(test_df.drop(columns=["id"]), columns=cat_cols) # Drop ID from
49 test set
50
51 # Align train/test features
52 X_aligned, test_aligned = X_enc.align(test_enc, join="left", axis=1, fill_value=0)
53
54 # Scale numeric features
55 scaler = StandardScaler()
56 X_scaled = scaler.fit_transform(X_aligned)
57 test_scaled = scaler.transform(test_aligned)
58
59 # === 4           Optuna Hyperparameter Tuning ===
60 def objective(trial):
61     params = {
62         "n_estimators": trial.suggest_int("n_estimators", 500, 2500),
63         "max_depth": trial.suggest_int("max_depth", 3, 10),
64         "learning_rate": trial.suggest_float("learning_rate", 0.01, 0.3, log=True),
65         "subsample": trial.suggest_float("subsample", 0.5, 1.0),
```

```

64     "colsample_bytree": trial.suggest_float("colsample_bytree", 0.5, 1.0),
65     "gamma": trial.suggest_float("gamma", 0, 5),
66     "min_child_weight": trial.suggest_int("min_child_weight", 1, 10),
67     "reg_alpha": trial.suggest_float("reg_alpha", 0.0, 5.0),
68     "reg_lambda": trial.suggest_float("reg_lambda", 0.0, 5.0),
69     "objective": "multi:softprob",
70     "num_class": len(le.classes_),
71     "eval_metric": "mlogloss",
72     "n_jobs": -1,
73     "tree_method": "hist",
74     "random_state": 42,
75     "verbosity": 0
76 }
77
78 model = XGBClassifier(**params)
79 cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
80 scores = cross_val_score(model, X_scaled, y_enc, cv=cv, scoring="accuracy", n_jobs=-1)
81 return scores.mean()
82
83 print("    Starting Optuna hyperparameter tuning...")
84 study = optuna.create_study(direction="maximize")
85 study.optimize(objective, n_trials=35, timeout=3600)
86
87 print("\n    Best Parameters Found:")
88 print(study.best_params)
89 print(f"Best CV Accuracy: {study.best_value:.4f}")
90
91 # === 5          Final Training with Best Params + K-Fold ===
92 best_params = study.best_params
93 best_params.update({
94     "objective": "multi:softprob",
95     "num_class": len(le.classes_),
96     "eval_metric": "mlogloss",
97     "n_jobs": -1,
98     "tree_method": "hist",
99     "random_state": 42,
100 })
101
102 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=42)
103 oof_preds = np.zeros(len(X_scaled))
104 test_preds_agg = np.zeros((len(test_scaled), len(le.classes_)))
105
106 print("\n    Starting Final 10-Fold Training...")
107 for fold, (train_idx, val_idx) in enumerate(skf.split(X_scaled, y_enc)):
108     model = XGBClassifier(**best_params)
109     model.fit(X_scaled[train_idx], y_enc[train_idx])
110     preds_val = np.argmax(model.predict_proba(X_scaled[val_idx]), axis=1)
111     oof_preds[val_idx] = preds_val
112     test_preds_agg += model.predict_proba(test_scaled) / skf.n_splits
113     acc = accuracy_score(y_enc[val_idx], preds_val)
114     print(f"Fold {fold+1} Accuracy: {acc:.4f}")
115
116 final_acc = accuracy_score(y_enc, oof_preds)
117 print(f"\n    Final Mean CV Accuracy: {final_acc:.4f}")
118
119 # === 6          Final Predictions + Save to CSV ===
120 final_test_preds = np.argmax(test_preds_agg, axis=1)
121 pred_labels = le.inverse_transform(final_test_preds.astype(int))
122
123 submission = pd.DataFrame({
124     "id": test_ids,
125     "TARGET": pred_labels
126 })
127
128 save_path = "/mnt/drive/MyDrive/Obesity Dataset/submission_optuna_xgb.csv"
129 submission.to_csv(save_path, index=False)
130 print(f"\n    Saved final predictions to: {save_path}")

```



```
131 print(submission.head())
```

Listing 2: ML Pipeline Script 1 (XGBoost + Optuna)

C Modeling Script 2: Random Forest + GridSearchCV

```
1
2 # =====
3 # Obesity Classification      Random Forest + GridSearchCV
4 # =====
5
6 import pandas as pd
7 import numpy as np
8 from sklearn.preprocessing import LabelEncoder, StandardScaler
9 from sklearn.model_selection import train_test_split, GridSearchCV
10 from sklearn.metrics import accuracy_score, classification_report
11 from sklearn.ensemble import RandomForestClassifier
12 from google.colab import drive
13 import warnings
14 warnings.filterwarnings("ignore")
15
16 # === 1 Mount Drive & Load Data ===
17 drive.mount('/mnt/drive')
18 train_df = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/train.csv')
19 test_df = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/test.csv')
20 sample_df = pd.read_csv('/mnt/drive/MyDrive/Obesity Dataset/sample_submission.csv')
21
22 print(f" Data Loaded: Train shape: {train_df.shape}, Test shape: {test_df.shape}")
23
24 # === 2 Preserve Test IDs ===
25 test_ids = test_df["id"].copy()
26
27 # === 3 Prepare Train/Test Data ===
28 TARGET = "WeightCategory"
29 X = train_df.drop(columns=["id", TARGET])
30 y = train_df[TARGET]
31
32 # Encode target
33 le = LabelEncoder()
34 y_enc = le.fit_transform(y)
35
36 # One-hot encode categorical features
37 cat_cols = X.select_dtypes(include=["object", "category"]).columns.tolist()
38 X_enc = pd.get_dummies(X, columns=cat_cols)
39 test_enc = pd.get_dummies(test_df.drop(columns=["id"]), columns=cat_cols)
40
41 # Align train/test features
42 X_aligned, test_aligned = X_enc.align(test_enc, join="left", axis=1, fill_value=0)
43
44 # Scale features
45 scaler = StandardScaler()
46 X_scaled = scaler.fit_transform(X_aligned)
47 test_scaled = scaler.transform(test_aligned)
48
49 # === 4 Train/Validation Split ===
50 X_tr, X_val, y_tr, y_val = train_test_split(
51     X_scaled, y_enc, test_size=0.2, random_state=42, stratify=y_enc
52 )
53
54 # === 5 Define Random Forest + Expanded GridSearchCV ===
55 rf = RandomForestClassifier(random_state=42, n_jobs=-1)
56
57 param_grid = {
58     "n_estimators": [200, 400],
59     "max_depth": [10, 15, None],
60     "min_samples_split": [2, 5],
61     "min_samples_leaf": [1, 2],
62     "max_features": ["sqrt", "log2"],
63     "bootstrap": [True],
64     "criterion": ["gini"],
65     "class_weight": [None, "balanced"]
66 }
```

```

66 }
67
68
69 grid_search = GridSearchCV(
70     estimator=rf,
71     param_grid=param_grid,
72     cv=5,                      # 5-fold cross validation
73     scoring="accuracy",
74     n_jobs=-1,
75     verbose=2
76 )
77
78 print("\n Running Expanded GridSearchCV...")
79 grid_search.fit(X_tr, y_tr)
80
81 print("\n Best parameters found:")
82 print(grid_search.best_params_)
83 print(f" Best CV Accuracy: {grid_search.best_score_:.4f}")
84
85 # Use best model
86 best_rf = grid_search.best_estimator_
87
88 # === 6 Evaluate on Validation Set ===
89 val_pred = best_rf.predict(X_val)
90 print(f"\n Validation Accuracy: {accuracy_score(y_val, val_pred):.4f}")
91 print("\n Classification Report:")
92 print(classification_report(y_val, val_pred, target_names=le.classes_))
93
94 # === 7 Retrain on Full Data ===
95 best_rf.fit(X_scaled, y_enc)
96
97 # === 8 Predict on Test Set ===
98 test_pred = best_rf.predict(test_scaled)
99 pred_labels = le.inverse_transform(test_pred)
100
101 # === 9 Build & Save Submission ===
102 submission = pd.DataFrame({
103     "id": test_ids,
104     TARGET: pred_labels
105 })
106
107 # Reorder columns to match sample submission if required
108 if set(sample_df["id"]) == set(submission["id"]):
109     submission = sample_df[["id"]].merge(submission, on="id", how="left")
110
111 save_path = "/mnt/drive/MyDrive/Obesity Dataset/submission_rf_grid_expanded.csv"
112 submission.to_csv(save_path, index=False)
113 print(f"\n Saved submission file to: {save_path}")
114 print(submission.head())

```

Listing 3: ML Pipeline Script 2 (Random Forest + GridSearchCV)