# Obesity Prediction Using Machine Learning Methods

Nidhi Mankala | Tris Marie Joe | Naveen Kumar N | SCOPE

## Introduction

In the current times, obesity is usually predicted using body mass index, but this diagnosis might not be always accurate. There are various other factors that can cause obesity. In our project we take into consideration all these external correlations to make predictions.

## Motivation

It is important to try various methods in public health sector and so we believe that incorporating Machine Learning will help to improve predictions.

## SCOPE of the Project

Analysis of different machine learning methods: Support Vector Machine, XGBoost, Adaboost, Gradient Boosting, Random forest classifier, category boosting and light gradient boosting. The models are trained on the dataset and is then evaluated using performance metrics such as accuracy, recall, precision score, ROC curve.

## Methodology

This work seeks to predict obesity using data from the 'Estimation of Obesity Levels based on Eating Habits and Physical Condition dataset' from the UCI Machine Learning Repository. There are a total of 2111 instances of data collected, with 17 properties.
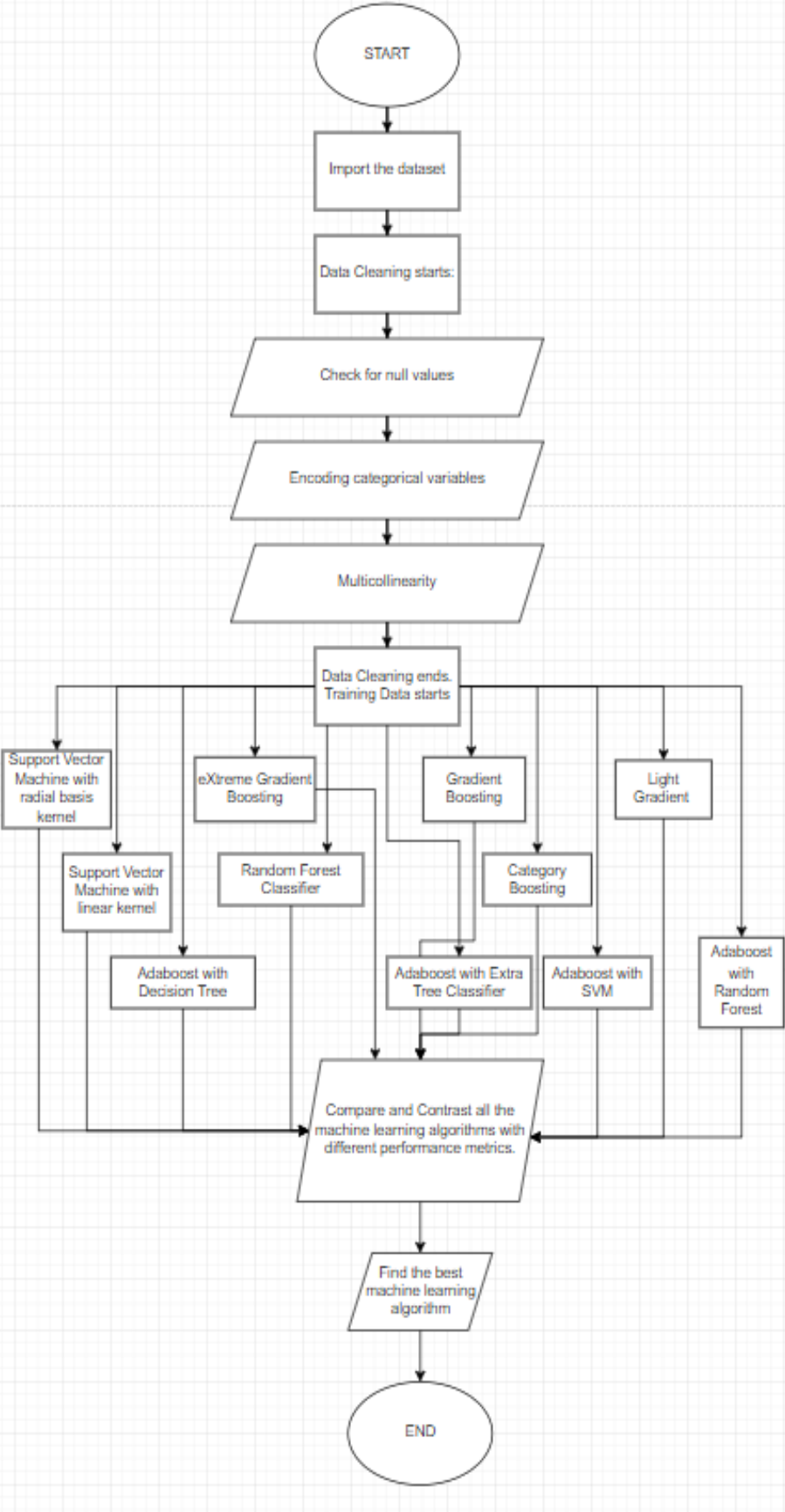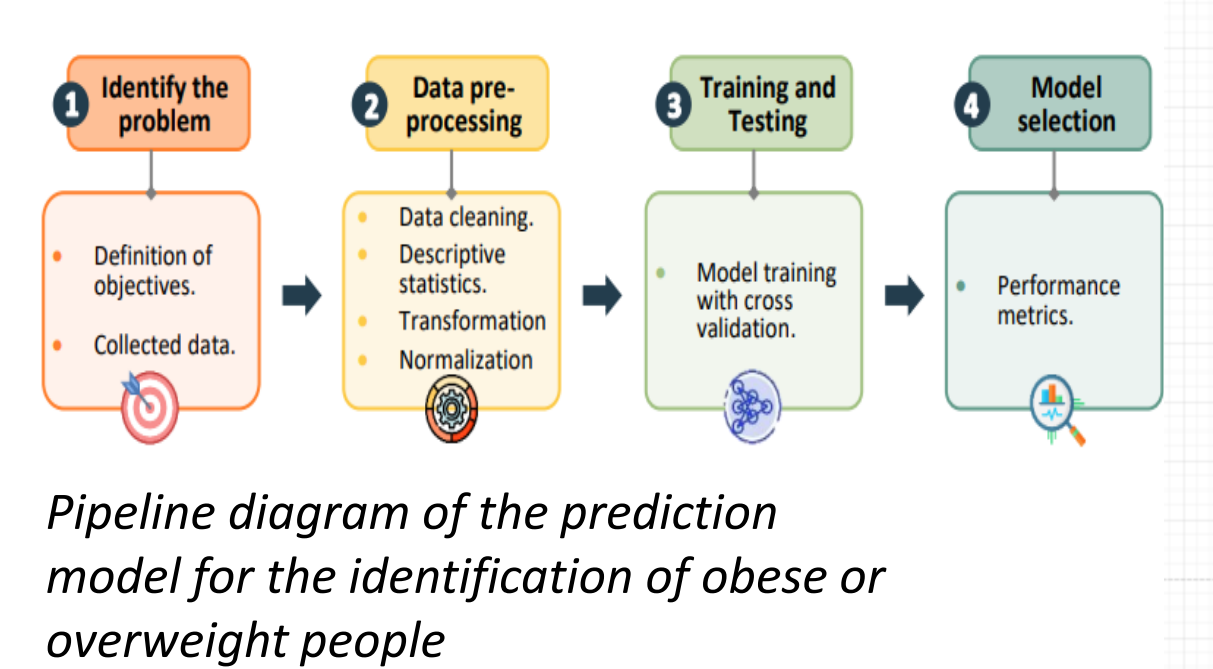
*Features of dietary habits*

| | | |
|---|---|---|
| 6 | c_ECFF | Do you eat high caloric food frequently? |
| 7 | c_EVM | Do you usually eat vegetables in your meals? |
| 8 | c_MMHD | How many main meals do you have daily? |
| 9 | c_EFBM | Do you eat any food between meals? |
| 10 | c_SMOKE | Do you smoke? |
| 11 | c_WDRD | How much water do you drink daily? |
| 12 | c_DRAL | How often do you drink alcohol? |

*List of esoteric features of the dataset that give a more accurate prediction*

*Physical condition features*

| | | |
|---|---|---|
| 13 | c_MCED | Do you monitor the calories you eat daily? |
| 14 | c_HPHA | How often do you have physical activity? |
| 15 | c_TTEC | How long do you use technological devices? |
| 16 | c_TRANSP | What type of transportation do you usually use? |

We start the project by importing our dataset, performing some methods for data cleaning to make our dataset ready. We then tried various different machine learning algorithms on this dataset and then compare the results of these algorithms with different performance metrics.
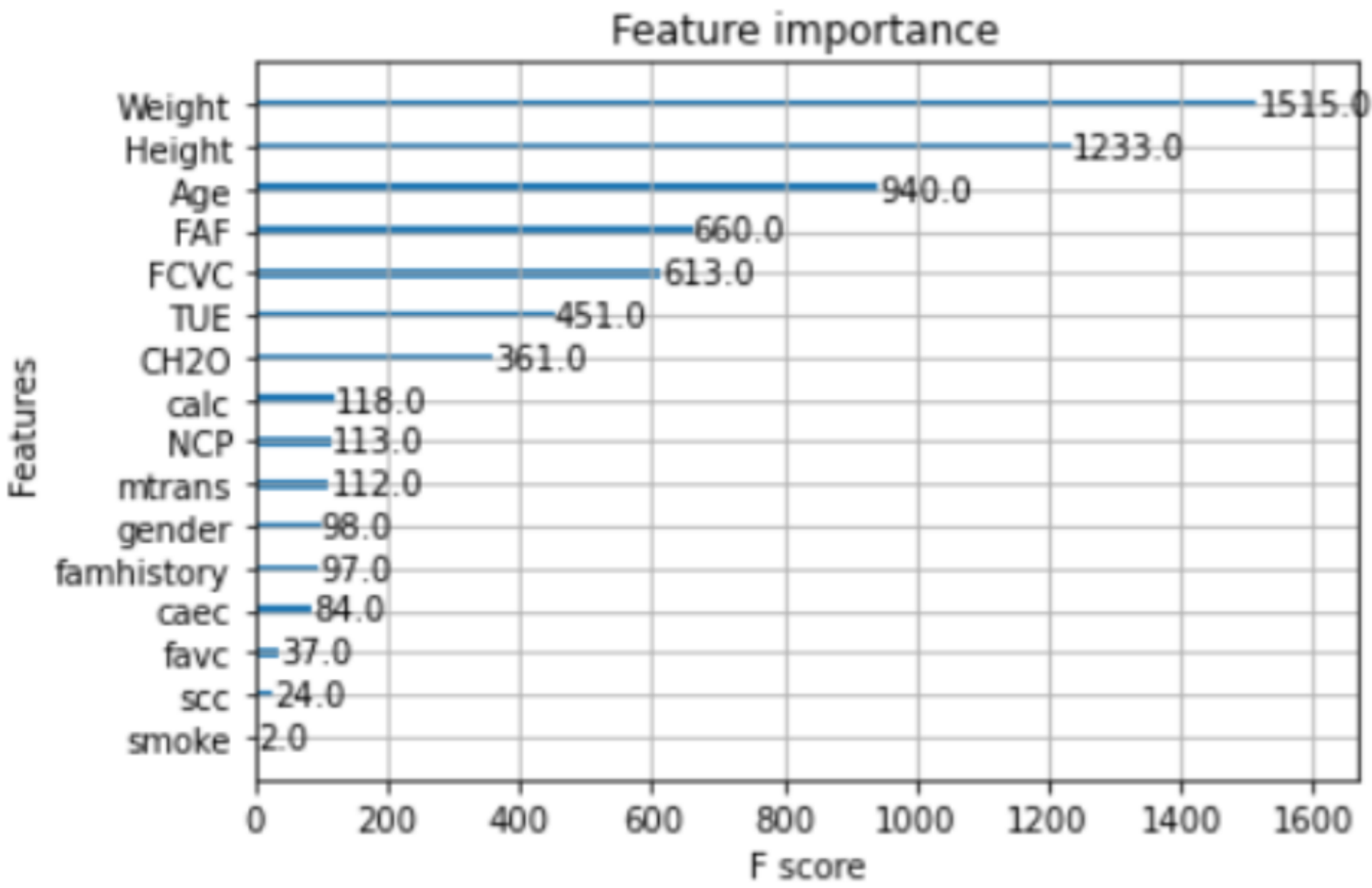


*Pipeline diagram of the prediction model for the identification of obese or overweight people*



*Architecture of Project*

## Results



*Correlation matrix of all the features in the dataset*



*Feature importance graph shows the feature which has the highest correlation to obesity.*

*Comparative analysis of different performance metrics on all the algorithms*

| ML MODELS | accuracy | precision | recall | f1 |
|---|---|---|---|---|
| SVM with radial basis kernel | 77.3 | 76.7 | 77.3 | 76.8 |
| SVM with linear kernel | 96.9 | 96.9 | 96.7 | 96.8 |
| eXtreme Gradient Boosting | 97.6 | 97.6 | 97.55 | 97.57 |
| Adaboost with Decision Tree | 93.1 | 92.9 | 92.8 | 92.8 |
| Adaboost with Random Forest | 96.2 | 96.4 | 96.1 | 96.1 |
| Random Forest Classifier | 95.5 | 95.6 | 95.3 | 95.3 |
| Gradient Boosting | 95.9 | 95.8 | 95.83 | 95.83 |
| Adaboost with Extra Tree classifier | 93.6 | 93.6 | 93.3 | 93.4 |
| Adaboost with SVM | 82.2 | 82.5 | 82.09 | 82.1 |
| Category Boosting (CatBoost) | 97.8 | 97.9 | 97.8 | 97.8 |
| Light Gradient Boosting Method | 99.05 | 99.04 | 99.09 | 99.06 |

## Conclusion

From our comparative analysis we see that the Light Gradient Boosting Method which gives an accuracy of 99.09% is the best ML algorithm for this prediction based model.

We have also optimized the results of other machine learning models to give a accurate prediction.

This product will help future developers use this prediction model to create product that can help people check their health digitally, hence contributing to smoother economic transactions.

## References

1. & Sánchez Hernández, A. B. (2019). Obesity level estimation software based on decision trees.
2. Cervantes, R. C., & Palacio, U. M. (2020). Estimation of obesity levels based on computational intelligence. *Informatics in Medicine Unlocked, 21,* 100472..