

Data-Driven Analytics for Obesity Management and Business Strategy

Author: Zhetan Zhang, Mingqi Yang

Team: Team ZY

1. Executive Summary

1.1 Decisions to be impacted

- Does this person have obesity?
- How to prevent obesity by changing lifestyle?
- What action should companies do to help obese people?
- What are the most important features that cause obesity?

1.2 Business value

- Nutrition companies and medical companies can design products that are more acceptable and effective for obese population based on their dietary habits
- Rehabilitation centers can provide more personalized services for obesity in order to broaden their prospective customers
- Internet companies can design mobile applications that can help potentially obese people to monitor their weights.

1.3 Data assets

- Estimation of obesity levels based on eating habits and physical condition . (2019). UCI Machine Learning Repository. <https://doi.org/10.24432/C5H31Z>.

This data set included an abundant set of variables, we can generate a wide-ranging view about how various factors collectively contribute to obesity. For example, by analyzing the features related to individuals' eating habits and their corresponding obesity levels, we can offer guidance to nutrition companies on developing more effective slimming products that are not only appealing to obese population but also effective in reducing weight. Besides, examining physical condition features also provide insights on how to suggest infrastructure companies such as public transit services or medical rehabilitation centers on optimizing convenience or construct appropriate exercise regimens for obese populations.

2. Data Preprocessing

2.1 Data Description

The dataset we are using for this project is called "Estimation of obesity levels based on eating habits and physical condition," sourced from the UCI Machine Learning Repository. The dataset essentially serves as classification tasks containing 16 features and 2111 samples, with the goal of predicting obesity levels in different individuals. Among 16 features, 6 of them contain nominal data (e.g. Gender, Transportation Used), 2 of them contain ordinal data (e.g.

Consumption of food between meals), 1 of them contains interval data, which is Age, and 7 of them contain ratio data. (e.g. Height, Weight)

To visualize the data, we separate it into two parts—numerical data and categorical data. For numerical data, we use bar plots to show the cumulative distribution of the data set. We have two bar plots which show the distribution of the raw data and the data after normalizing. (See figure 1 and figure 2)

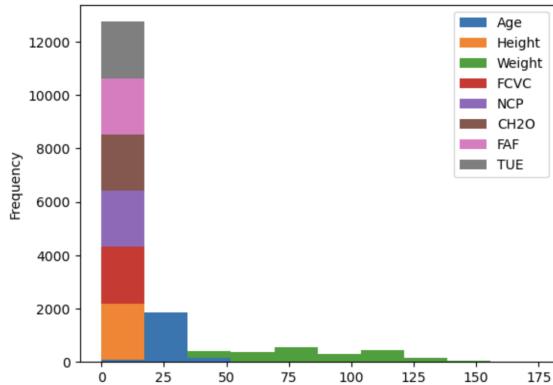


Figure 1: Raw data

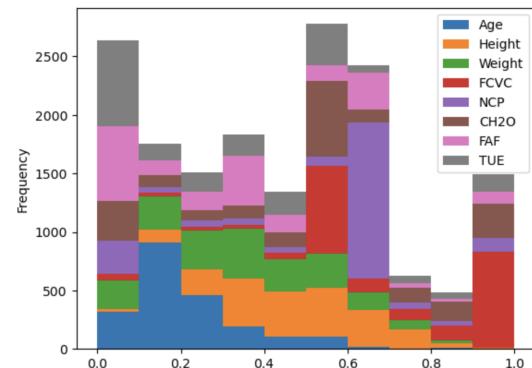


Figure 2: Normalized data

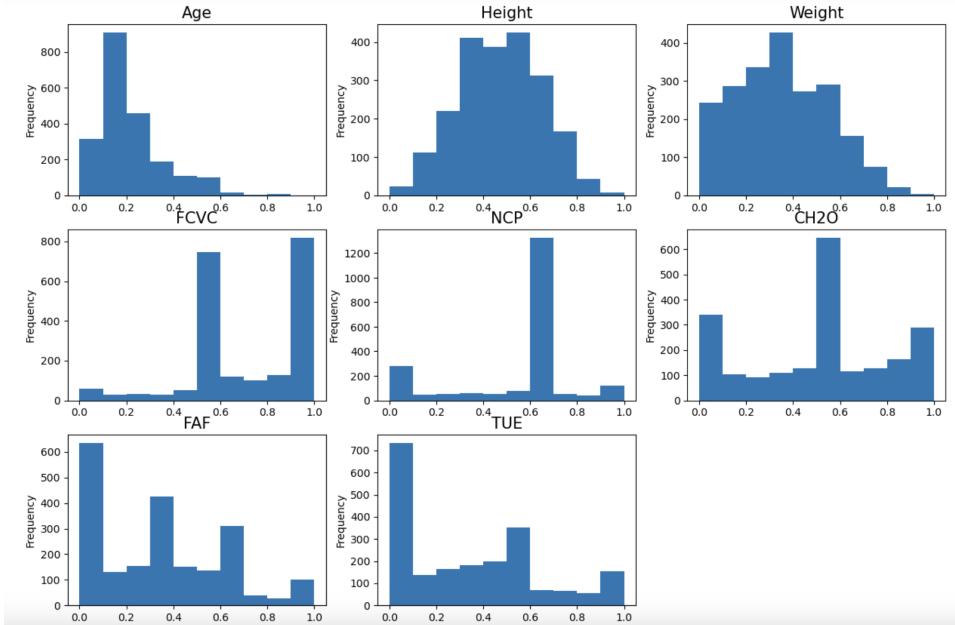


Figure 3: Data distribution

After normalizing, we can see that the data are all shrink into the range [0,1] and have better distribution than the raw data. Then, we can use correlation matrix to see if there any features have high correlation.(See figure 4)

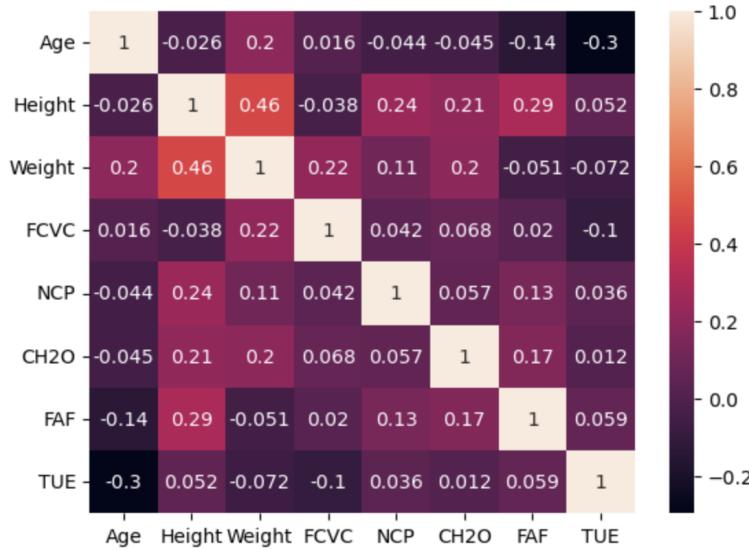


Figure 4: Correlation Heat map

From the correlation heat map, we can see there are no pairs of numeric features that have a correlation greater or equal to 0.8 which means that we can assume each numeric feature is independent of each other.

For categorical data, we used pie charts to show the distribution.(See figure 5)

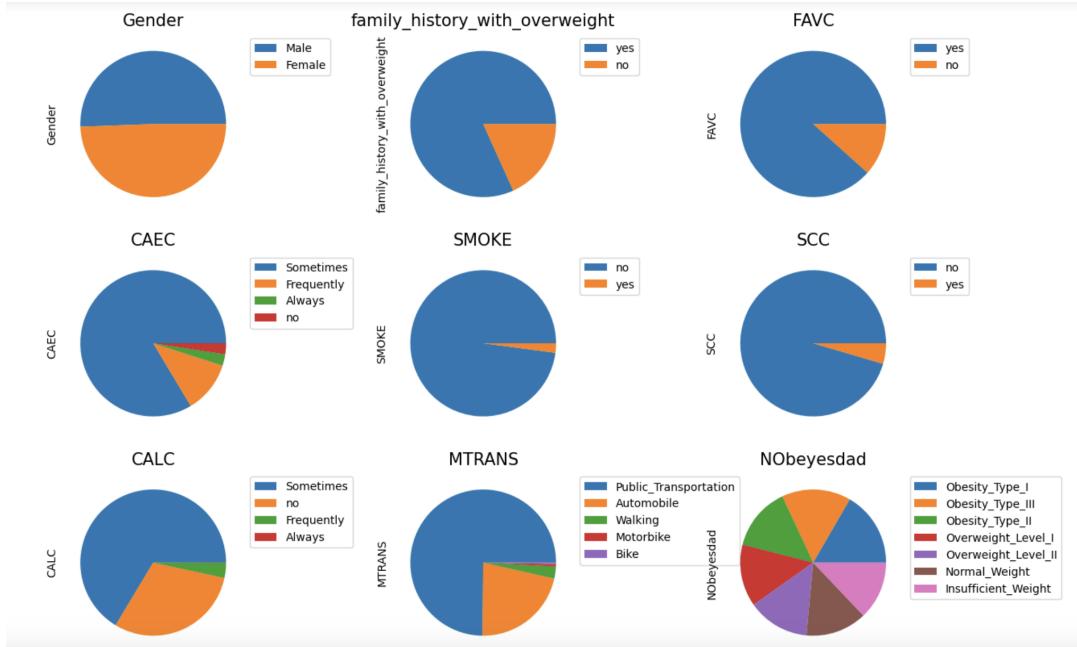


Figure 5: Categorical data Pie charts

2.2 Data Cleaning

We used one-hot-encoding to encode categorical data to transform categorical data into numeric data. For each unique category in the categorical variable, one-hot-encoding will create a new binary column. In the new binary columns, you assign a value of 1 if the original data point belongs to the category represented by that column, and 0 if it does not. These binary columns form a new feature matrix, where each row represents a data point, and each column represents a category from the original categorical variable.

2.3 Outlier Detection

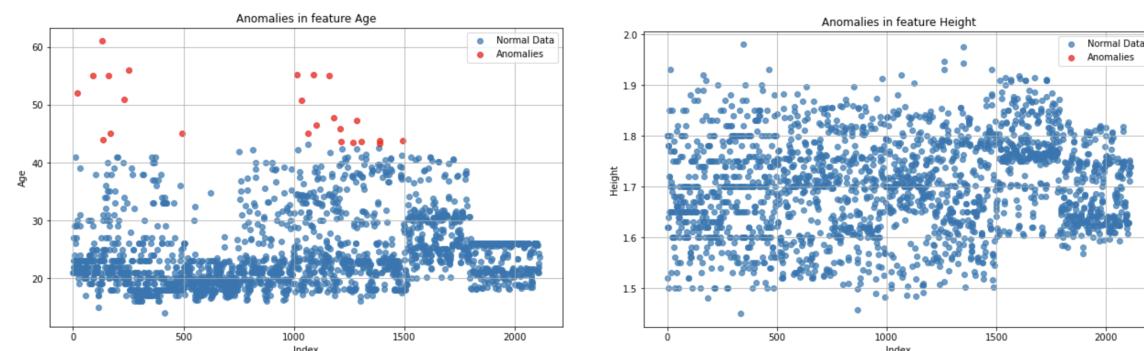
For identifying outliers in our dataset, we primarily employed different statistical methods on both numerical and categorical features. Specifically, for each numerical attribute, we calculate its corresponding standard deviation and mean value, which has shown below.

	Age	Height	Weight	FCVC	NCP	CH2O	FAF	TUE
count	2111.000000	2111.000000	2111.000000	2111.000000	2111.000000	2111.000000	2111.000000	2111.000000
mean	24.312600	1.701677	86.586058	2.419043	2.685628	2.008011	1.010298	0.657866
std	6.345968	0.093305	26.191172	0.533927	0.778039	0.612953	0.850592	0.608927
min	14.000000	1.450000	39.000000	1.000000	1.000000	1.000000	0.000000	0.000000
25%	19.947192	1.630000	65.473343	2.000000	2.658738	1.584812	0.124505	0.000000
50%	22.777890	1.700499	83.000000	2.385502	3.000000	2.000000	1.000000	0.625350
75%	26.000000	1.768464	107.430682	3.000000	3.000000	2.477420	1.666678	1.000000
max	61.000000	1.980000	173.000000	3.000000	4.000000	3.000000	3.000000	2.000000

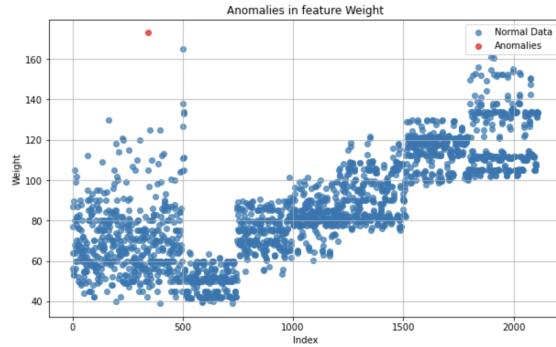
Then, we consider the feature values that are 3 standard deviations from the mean as anomalies. The graphs below shows the results of outlier detection on numeric

Feature Age:
Number of normal data points: 2087
Number of anomalies: 24
Percentage of anomalies: 1.149976042165788%

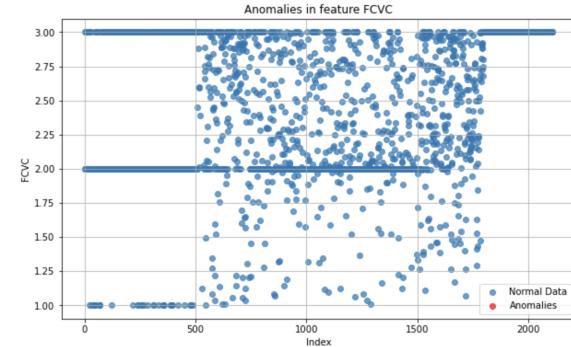
Feature Height:
Number of normal data points: 2111
Number of anomalies: 0
Percentage of anomalies: 0.0%



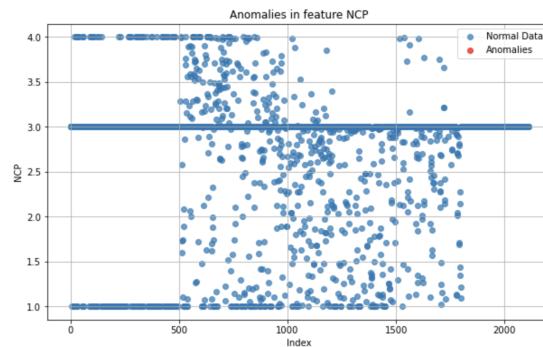
Feature Weight:
 Number of normal data points: 2110
 Number of anomalies: 1
 Percentage of anomalies: 0.047393364928909956%



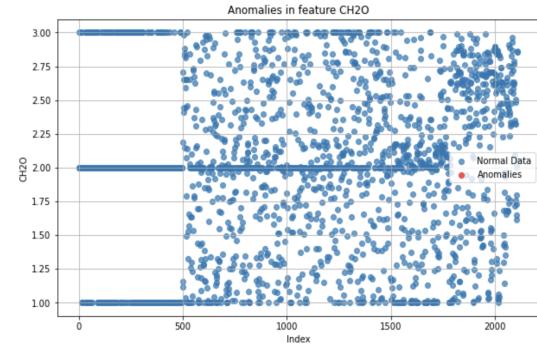
Feature FCVC:
 Number of normal data points: 2111
 Number of anomalies: 0
 Percentage of anomalies: 0.0%



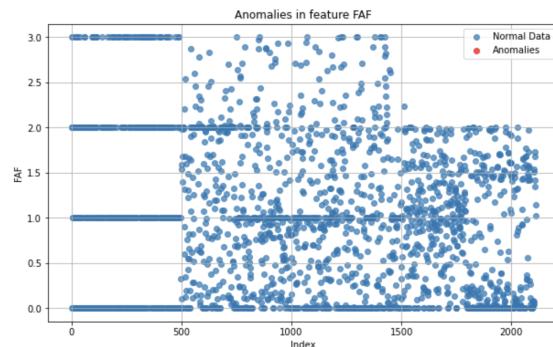
Feature NCP:
 Number of normal data points: 2111
 Number of anomalies: 0
 Percentage of anomalies: 0.0%



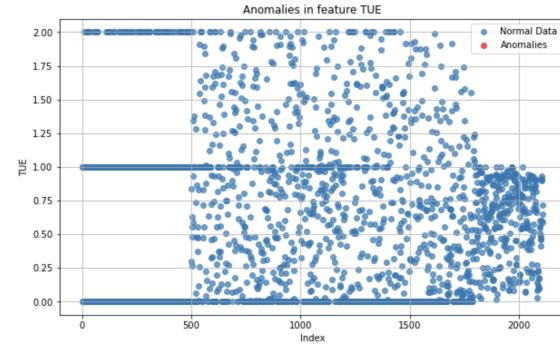
Feature CH20:
 Number of normal data points: 2111
 Number of anomalies: 0
 Percentage of anomalies: 0.0%



Feature FAF:
 Number of normal data points: 2111
 Number of anomalies: 0
 Percentage of anomalies: 0.0%



Feature TUE:
 Number of normal data points: 2111
 Number of anomalies: 0
 Percentage of anomalies: 0.0%



Based on the graphs above, we can see that for most of the features, there is no outlier within the dataset. However, when examining the “age” feature, we observed that samples with ages above 40 are outliers; when examining the “weight” feature, we observe only one sample with weight over 160 kg is outlier. To address “age” outliers, we intend to initially apply models like logistic regression to assess their effect on the models’ accuracy. For the outlier in “weight” feature, we decide to ignore it because it is unlikely to significantly impact the general performance of the model.

For each categorical feature, we have chosen to classify feature values that appears fewer than 5 times as the outliers. Given that our dataset comprises around 2000 samples, we believe this outlier definition is justifiable. The graph below shows the results of outlier detection on categorical dataset:

```
Anomalies in Gender:  
Series([], Name: Gender, dtype: int64)  
  
Anomalies in family_history_with_overweight:  
Series([], Name: family_history_with_overweight, dtype: int64)  
  
Anomalies in FAVC:  
Series([], Name: FAVC, dtype: int64)  
  
Anomalies in CAEC:  
Series([], Name: CAEC, dtype: int64)  
  
Anomalies in SMOKE:  
Series([], Name: SMOKE, dtype: int64)  
  
Anomalies in SCC:  
Series([], Name: SCC, dtype: int64)  
  
Anomalies in CALC:  
Always 1  
Name: CALC, dtype: int64  
  
Anomalies in MTRANS:  
Series([], Name: MTRANS, dtype: int64)  
  
Anomalies in NObeyesdad:  
Series([], Name: NObeyesdad, dtype: int64)
```

To be specific, we try to create panda series for the outliers within each feature. Based on the graph, you can see that for most of the features, no outlier exists. However, for the feature ‘Consumption of alcohol (CALC)’, we identified a sample with the value "always" as an anomaly. Similar as before, we decide to ignore this single outlier because it will not have a significant effect on the model’s performance.

We then used LOF(Local outlier factories) to detect outliers. LOF is a density-based method that identifies outliers by comparing the local density deviation of a data point with respect to its neighbors. The basic idea is that outliers are points that have a significantly lower local density compared to their neighbors. We used the default parameter of LOF in sklearn which $n_neighbor = 20$. By using LOF, we detected 100 samples as outliers which is $100/2111 = 4.7\%$ data. We will use the data set which delete these outliers and compare with the performance of the original data.

3. Model Updates

3.1 Logistic Regression

Our first purpose is to use some features to predict whether a person has obesity or not. This problem should be a binary classification problem and we need to separate the target feature

“NObeyesdad” into two parts. One part should be “obesity” including all the obesity types as 1(positive class) and “no obesity” including normal weight and insufficient weight as 0(negative class). The model we choose to apply first is logistic regression since logistic regression is a powerful model in binary classification. In addition, by using logistic regression we can see the importance of each feature. This can help us figure out which factors affect obese people the most and help them change their lifestyle. Before applying logistic regression, we split the dataset into training(70%) and testing(30%) set by using train_test_split in the sklearn package. By applying logistic regression, we have the following testing result.

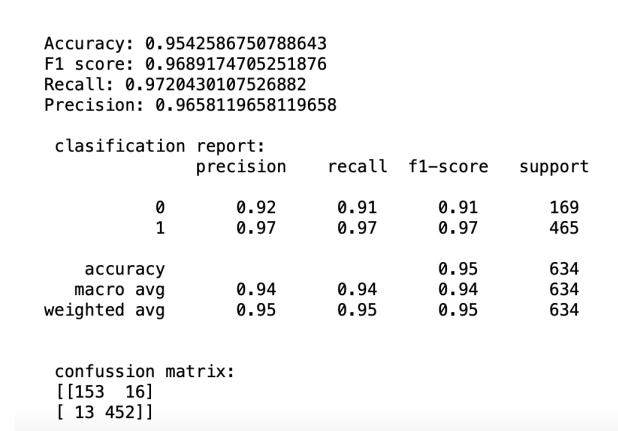


Figure 6: Evaluation matrix

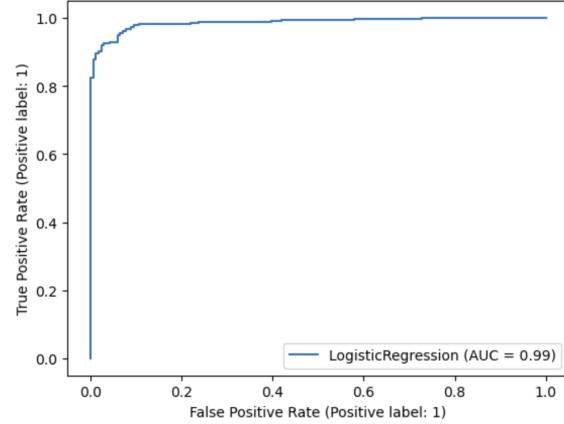


Figure 7: ROC

We can see the accuracy of this model is around 95% and the area under the curve is 0.99 which is really good. Besides this, we can see the weight of each feature. By checking the absolute value of each weight, we have the largest top 3 absolute weight features: “Weight”, “Age”, “Height” which make sense since these are the main factors to determine if a person is overweight or not. See figure 8 as the whole weight list of features.

Weight	11.638546	CALC_Frequently	0.498513
Age	2.771782	CH2O	0.435136
Height	-2.508798	SMOKE_no	0.422475
CAEC_no	1.007944	SMOKE_yes	-0.422094
CAEC_Frequently	-0.990376	SCC_yes	0.302898
CAEC_Always	-0.914366	SCC_no	-0.302517
CAEC_Sometimes	0.897178	CALC_no	-0.297422
FAF	-0.833509	FCVC	-0.289440
MTRANS_Public_Transportation	0.794503	TUE	-0.217558
NCP	-0.682140	MTRANS_Motorbike	-0.199483
MTRANS_Bike	-0.628362	FAVC_yes	0.175246
MTRANS_Automobile	0.609239	FAVC_no	-0.174865
family_history_with_overweight_yes	0.578395	CALC_Sometimes	-0.109209
family_history_with_overweight_no	-0.578014	CALC_Always	-0.091502
MTRANS_Walking	-0.575516	Gender_Female	0.068078
		Gender_Male	-0.067697

Figure 8: Feature weight

We can also see that for “CAEC_no” and “CAEC_Sometimes”, they have positive weight. “CAEC_Frequently” and “CAEC_Always” have negative weights. CAEC is the Consumption of food between meals. The positive weight means this term will affect more on predicting the positive result(obesity) and negative weights will affect more on predicting the negative result(not obesity). By this model, we can conclude that more frequent consumption of food between meals will more likely lead to obesity. By using this pattern, we also can conclude the following conditions will likely lead to obesity. Low Physical activity frequency(FAF), less Number of main meals(NCP), use public transportation or automobile, family history with overweight, frequently Consumption of alcohol(CALC), more Consumption of water daily(CH2O), do not smoke, will do Calories consumption monitoring(SCC), Frequent consumption of high caloric food(FAVC).

We also applied Logistic regression on the data cleaned by LOF and we got the following results.

```

Accuracy: 0.945364238410596
F1 score: 0.96353591160221
Recall: 0.9819819819819819
Precision: 0.9457700650759219

classification report:
precision    recall   f1-score   support
0            0.94      0.84      0.89      160
1            0.95      0.98      0.96      444

accuracy           0.95      604
macro avg       0.94      0.91      0.93      604
weighted avg     0.95      0.95      0.94      604

confusion matrix:
[[135 25]
 [ 8 436]]

```

We can find there is no large difference of evaluate value between the original model.

We also used cross validation on logistic regression to make sure there did not exist overfitting and here is the result we got from the 5-fold cross validation.

```

{'mean_fit_time': array([0.01663008]),
 'std_fit_time': array([0.01131471]),
 'mean_score_time': array([0.00049195]),
 'std_score_time': array([0.00022771]),
 'params': [{}],
 'split0_test_score': array([0.95945946]),
 'split1_test_score': array([0.94256757]),
 'split2_test_score': array([0.92542373]),
 'split3_test_score': array([0.94915254]),
 'split4_test_score': array([0.93559322]),
 'mean_test_score': array([0.9424393]),
 'std_test_score': array([0.01158593]),
 'rank_test_score': array([1], dtype=int32),
 'split0_train_score': array([0.95427604]),
 'split1_train_score': array([0.9491956]),
 'split2_train_score': array([0.94923858]),
 'split3_train_score': array([0.94670051]),
 'split4_train_score': array([0.9500846]),
 'mean_train_score': array([0.94989906]),
 'std_train_score': array([0.00246391])}

```

We can find there is no large difference between test score and train score; therefore, we can conclude that our model has no overfitting.

By the result we got before, we can find terms “weight”, “height”, and “age” have much larger weight than other terms. Therefore, we want to drop these three terms to build a new model and this model can be interpreted as predicting obesity type only by using people’s lifestyle. This model can solve some decisions to be impacted we mentioned at the beginning of this report. We got the following result.

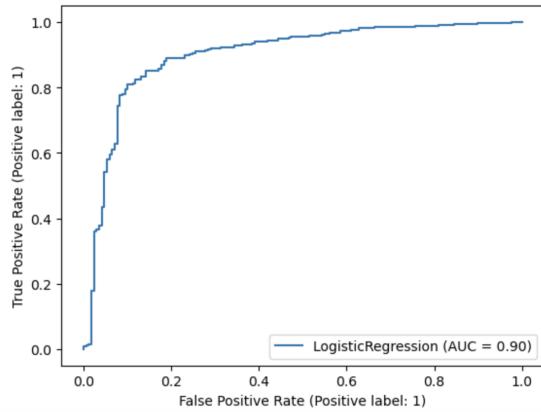
```
Accuracy: 0.8580441640378549
F1 score: 0.9050632911392404
Recall: 0.9225806451612903
Precision: 0.8881987577639752

classification report:
precision    recall   f1-score   support
0            0.76     0.68     0.72      169
1            0.89     0.92     0.91      465

accuracy       0.82
macro avg      0.82
weighted avg   0.85

confusion matrix:
[[115  54]
 [ 36 429]]
```

Evaluation matrix



ROC

From the result, we can find there is an obvious decrease of the precision between the original model and the new model. Therefore, we might need to consider some complexer models to make predictions by dropping “weight”, “height”, and “age”.

3.2 K-Nearest Neighbors

The second model we used is Knn. Knn is a simple but widely used algorithm in machine learning, especially in classification problems. There exists an important hyperparameter in Knn which is the “n_neighbors”. We used grid search to choose the best “n_neighbors” and used cross validation to make sure there did not exist heavy overfitting. We build two different models of Knn, one is the same as Logistic regression which separates response value to Obese or not which is binary classification, another is to use the original response value which has 7 different values. We will see the performance of both of these models.

For both binary classification and multivariate classification , we got the best “n_neighbors = 1”. We got the following results,

```

Accuracy: 0.7870662460567823
F1 score: 0.7805979049752149
Recall: 0.7831841848685303
Precision: 0.7884924177220066

clasification report:
precision      recall   f1-score   support
                                         1       0.62     0.78     0.69     89
                                         2       0.66     0.84     0.74     85
                                         3       0.86     0.75     0.80     114
                                         4       0.92     0.94     0.93     85
                                         5       0.99     0.98     0.98     92
                                         6       0.59     0.47     0.52     77
                                         7       0.88     0.74     0.80     92

clasification report:
precision      recall   f1-score   support
                                         accuracy    0.79     0.79     634
                                         macro avg  0.79     0.79     634
                                         weighted avg 0.80     0.79     634
                                         0       0.86     0.86     0.86     169
                                         1       0.95     0.95     0.95     465

accuracy      0.92     0.92     0.92     634
macro avg    0.90     0.90     0.90     634
weighted avg 0.92     0.92     0.92     634

confussion matrix:
[[69  7  3  2  0  8  0]
 [11 71  1  0  1  1  0]
 [ 8 14 85  1  0  4  2]
 [ 2  2  1 80  0  0  0]
 [ 0  2  0  0 90  0  0]
 [14 11  7  2  0 36  7]
 [ 7  1  2  2  0 12 68]]

```

Binary classification

Multi-class classification

We can find the binary classification performance is worse than Logistic regression but have a pretty good performance overall. However, for multivariable classification, some variables have high precision like 5:Obesity_Type_III and some variables have low precision like 6:Normal_Weight. This phenomenon leads to a not good overall performance. We might not consider using Knn as the final model for multivariate classification.

We also apply Knn on cleaned data and we got the following results:

```

Accuracy: 0.8162251655629139
F1 score: 0.8035116189598093
Recall: 0.8071648884289305
Precision: 0.8043031549025633

clasification report:
precision      recall   f1-score   support
                                         1       0.74     0.71     0.72     76
                                         2       0.69     0.79     0.74     77
                                         3       0.81     0.82     0.82     100
                                         4       0.91     0.94     0.92     93
                                         5       0.99     1.00     0.99     98
                                         6       0.68     0.52     0.59     81
                                         7       0.81     0.87     0.84     79

clasification report:
precision      recall   f1-score   support
                                         accuracy    0.82     0.82     604
                                         macro avg  0.80     0.81     604
                                         weighted avg 0.81     0.82     604
                                         0       0.85     0.84     0.85     160
                                         1       0.94     0.95     0.95     444

accuracy      0.92     0.92     0.92     604
macro avg    0.90     0.89     0.90     604
weighted avg 0.92     0.92     0.92     604

confussion matrix:
[[54  7  5  0  0  7  3]
 [ 9 61  3  1  1  1  1]
 [ 0  8 82  5  0  4  1]
 [ 0  2  4 87  0  0  0]
 [ 0  0  0  0 98  0  0]
 [10 10  5  3  0 42 11]
 [ 0  0  2  0  0  8 69]]

```

Binary classification

Multi-class classification

By the result, we can see this is the same as Logistic regression, there did not exist some huge difference compared to the original model. For the multi-class classification, some classes have a slight improvement of precision.

3.3 Random Forest Tree

The third model we employed is the Random Forest Tree algorithm. Random Forest Tree is a learning-based approach that incorporates various decision trees to construct a “forest” and generate predictions. Specifically, each decision tree in the “forest” will make their own predictions based on a subset of data and features, and the “forest” will choose the most common outcomes among all trees to determine the final decision. Like KNN, Random Forest Tree can be used for both binary classification and multi-class classification. Compared to KNN, Random Forest Tree is more robust in handling the larger datasets and non-linear relationships. However, Random Forest Tree needs more computational resources and training time due to the construction of several decision trees.

To determine if an individual is obese and identify the specific obesity type they may have, we utilize the Random Forest Tree for both binary and multi-class classification tasks. The results for each task are shown in the figure below:

```
Accuracy: 0.9258675078864353
F1 score: 0.9234036905903121
Recall: 0.9257840540235982
Precision: 0.9267647472384187

classification report:
precision    recall   f1-score   support
1            0.95     0.78     0.85      89
2            0.86     0.93     0.89      85
3            0.99     0.92     0.95     114
4            1.00     1.00     1.00      85
5            1.00     1.00     1.00      92
6            0.74     0.91     0.81      77
7            0.96     0.95     0.95      92

Accuracy: 0.9700315457413249
F1 score: 0.9795479009687836
Recall: 0.978494623655914
Precision: 0.9806034482758621

clasification report:
precision    recall   f1-score   support
accuracy
0            0.94     0.94     169      634
1            0.98     0.98     465      634

accuracy      0.97     634
macro avg     0.96     634
weighted avg  0.97     634

accuracy
macro avg     0.96     634
weighted avg  0.97     634

confussion matrix:
[[ 69  4  0  0  0  16  0]
 [ 3  79  1  0  0  2  0]
 [ 1  6 105  0  0  2  0]
 [ 0  0  0  85  0  0  0]
 [ 0  0  0  0  92  0  0]
 [ 0  3  0  0  0  70  4]
 [ 0  0  0  0  0  5  87]]
```

(Binary Classification)

(Multi-class Classification)

For now, we have not done any hyperparameter tuning yet. In other words, we rely on the default settings of scikit-learn package’s Random Forest Classifier to make predictions. Compared to KNN and Logistic Regression, we can see that Random Forest Tree performs better in both classification tasks even with the default settings. However, we can still observe that Random Forest tree does not perform well in classifying “6: Normal Weight” category in the multi-class classification. In light of this, we need to dig into the samples that it wrongly classified to get more information.

We also perform the Random Forest Tree algorithm on the dataset that is free from outliers. The result for both classification task have been shown below:

```

Accuracy: 0.9420529801324503
F1 score: 0.9384146004837258
Recall: 0.9387776418489481
Precision: 0.9389362411613318

classification report:
precision    recall   f1-score   support
1            0.92    0.86     0.88      76
2            0.90    0.95     0.92      77
3            0.99    0.94     0.96     100
4            1.00    0.99     0.99      93
5            0.99    1.00     0.99      98
6            0.81    0.86     0.84      81
7            0.96    0.97     0.97      79

clasification report:
precision    recall   f1-score   support
0            0.96    0.94     0.95     160
1            0.98    0.99     0.98     444
accuracy          0.97    604
macro avg       0.97    0.96     0.97     604
weighted avg    0.97    0.97     0.97     604

accuracy          0.97    604
macro avg       0.97    0.96     0.97     604
weighted avg    0.97    0.97     0.97     604

confussion matrix:
[[65  2  1  0  0  8  0]
 [ 0 73  0  0  0  4  0]
 [ 2  2 94  0  0  2  0]
 [ 0  0  0 92  1  0  0]
 [ 0  0  0  0 98  0  0]
 [ 4  4  0  0  0 70  3]
 [ 0  0  0  0  0  2 77]]
```

(Binary Classification)

(Multi-class Classification)

Based on the graphs above, we can see that the outliers really have some negative effect on the performance of the Random Forest Tree and cleaning them up really improves the accuracy of the classifier.

Recall that Height, Weight and Age are three most important features for predicting obesity, we also drop these three features to see how the rest of the features interact with the prediction of obesity. The results of the prediction without Height, Weight and Age have shown in the graphs below:

```

Accuracy: 0.8079470198675497
F1 score: 0.8002110122444307
Recall: 0.800288931465978
Precision: 0.803746109049475

Accuracy: 0.9072847682119205
F1 score: 0.9370786516853933
Recall: 0.9391891891891891
Precision: 0.9349775784753364

clasification report:
precision    recall   f1-score   support
          1       0.77     0.67     0.72      76
          2       0.71     0.64     0.67      77
          3       0.79     0.75     0.77     100
          4       0.87     0.89     0.88      93
          5       0.99     1.00     0.99      98
          6       0.61     0.74     0.67      81
          7       0.88     0.91     0.89      79

          0       0.83     0.82     0.82     160
          1       0.93     0.94     0.94     444

accuracy           accuracy
macro avg       0.88       0.88       0.88      604
weighted avg    0.91     0.91     0.91     604

confusion matrix:
[[51  3  6  2  0 12  2]
 [5 49  8  5  0 10  0]
 [4  6 75  5  1  7  2]
 [1  4  2 83  0  3  0]
 [0  0  0  0 98  0  0]
 [5  7  3  0  0 60  6]
 [0  0  1  0  0  6 72]]

```

(Binary Classification)

(Multi-class Classification)

Obviously, we observe an apparent decrease in the performance of Random Forest Classifiers. However, we can notice that the performance of classifying “5: Obesity Type III” is still impressive with absence of Height, Weight, and Age. From this perspective, we may explore how the rest of the features contribute to the cause of Obesity Type III.

3.4 Machine Learning Workflow

The first thing we've done for the workflow is the data checking. Specifically, we acquired a dataset from UCI Machine Learning Repository, where the input features are elements of space \mathbf{X} and the target variables are elements of space \mathbf{Y} . Then, examined the dataset to see if there are any missing values and found no absent entries. If it does, we will perform a function f to handle these missing values,

$$\forall x \in X : x = \begin{cases} x, & \text{if } x \neq \text{null} \\ f(X), & \text{otherwise} \end{cases}$$

where x is the data point in \mathbf{X} and $f(X)$ is the function we apply to fulfill the missing values

Then, we perform descriptive analytics on the dataset by visualizing the distribution of each feature. In particular, we employ bar plots and histogram for numerical features and pie charts for categorical features. Besides, we also draw the correlation matrix to determine whether there are internal relationships among different variables.

After that, we applied one-hot-encoding on the categorical features to prepare them for our model analysis. Moreover, we normalized all the data values between 0 and 1, which is convenient for our potential linear model applications. Mathematically, suppose c represents each column of \mathbf{X} and \mathbf{D} represent the domain, we perform one-hot encoding like this:

$$\text{encode}(c) = \mathbf{D} \in [0, 1]$$

For now, we choose the logistic regression algorithm for our obesity level estimation. The reason is that our goal is to correctly classify the samples to be obese or not. Then, we randomly split the dataset into a training set and a test set, allocating 70% of the data for training and 30% of the data for testing. In MLM form, we can express our steps like this:

- The input space is $\mathbf{X} = \mathbb{R}^m$ and the output space is $\mathbf{Y} = \mathbb{R}$
- Since we do not embed prior knowledge for now, the error prior $P(\Theta) = I$
- The corresponding learning Morphism is $z = \Theta^T \mathbf{x}$
- The risk function is

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(p^{(i)}) + (1 - y^{(i)}) \log(1 - p^{(i)})]$$

$$\text{where } p = \frac{1}{1 + e^{-z}}$$

- Then, the optimal parameters are given by

$$\theta^* = \arg \min_{\theta} J(\theta)$$

and the MLM can be written as

$$ML_{LR} = \left(\mathbb{R}^m, \mathbb{R}, \theta^T \mathbf{x}, P(\theta) = 1, -y \log\left(\frac{1}{1 + e^{-\theta^T \mathbf{x}}}\right) - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\theta^T \mathbf{x}}}\right) \right)$$

4. Source Code

<https://github.com/Sonic-zzt/ESE-527-Project.git>

5. Next Steps

5.1 Hyperparameter Tuning

For now, we are just using the default setting and parameters of the Random Forest Tree algorithm. Although the algorithm shows generally remarkable performance, we still observe that some specific classes are not correctly classified compared to the others. In this case, we may need to perform the hyperparameter tuning to see whether we can further improve the performance of Random Forest Classifier.

5.2 Wrong Samples Analysis

Now we have seen that Random Forest Classifier shows better performance than the other two models, and we may choose it as our final model. So, it is the time for us to analyze the samples that are wrongly classified during the prediction. We should answer the question why these samples are different from others and determine whether they are potential outliers that are detrimental to the model's performance. Furthermore, we may think about a specific solution to deal with these samples when we make our final recommendation to the company.

5.3 Feature Engineering

Although Random Forest Tree shows impressive performance, we need to consider that we are only using a small curated dataset. In the real world, the datasets are often high-dimensional and large-sized. Therefore, we may want to perform feature engineering such as PCA on our original datasets to reduce the number of features and still maintain the variance of most features, which may make Random Forest Tree less complex and faster to construct. However, we still need to be careful about the loss of relevant information during the feature engineering process, which may lead to worse performance.