

Predicting BMI from Food Acquisition and Availability

Alan Liu

CIS 3200-02 Data Processing and Analytics
Department of Information Systems
California State University Los Angeles

Abstract: The United States Department of Agriculture conducted¹ the National Household Food Acquisition and Purchase Survey (FoodAPS) in the period between April 2012 and January 2013 [1]. As part of the dataset, the various types of food sources within different radii were counted and recorded per participating household, as well as individual health data, including their BMI.² Using Microsoft Azure ML Studio, the dataset was cleaned and used to train Bayesian Regression models in various configurations. The best performing model yielded a mean absolute error of 2.039. While fairly accurate relative to the mean BMI of the dataset,³ the low coefficient of determination⁴ suggests the frequency of various food sources by itself is not the optimal predictor of individual BMI.

1. Introduction

The goal of this project was to accurately predict an individual's BMI given the frequency of food sources near their household. In the survey data collected by the United States Department of Agriculture, the participants self-reported their food acquisition events for the period between April 2012 and January 2013. In addition to this, socioeconomic and health data was collected from the participants. The participants' households were geographically referenced to each and every known food source, including grocers, supermarkets, superstores, convenience stores, fast-food and non-fast food restaurants. It is the count of these food sources that will primarily be used to predict BMI.

2. Related Works and Background

This project attempts to use the non-medical factors of an individual's food environment to predict their BMI. As such, it uses primarily uses factors relating to the accessibility of food.

There are other non-medical studies in finding contributing factors to BMI. Daniel Fuller et. al. published their study in The American Journal of Clinical Nutrition in January 2013 about the relationship between transportation mode, distance to food stores, fruit and vegetable consumption, and BMI in low-income neighborhoods [2].

¹ USDA's Economic Research Service and Nutrition Service co-sponsored the survey.

² Body Mass Index, calculated as weight (kg) / height (m)².

³ The mean BMI in the dataset was approximately 25.

⁴ The Coefficient of Determination was 0.596.

The assumption was that using an automobile would convey easier access to food stores and would increase fruit and vegetable consumption, which would imply a healthier diet and lower BMI. Thus, those in low-income neighborhoods, being less likely to have access to an automobile, would have higher BMI because of presumably more difficult access to food stores and thus a healthy diet.

However, the study found no difference in fruit and vegetable consumption between different modes of transportation and found instead that those using public transportation had significantly lower BMI than compared to those who use an automobile [2]. This suggests that there are factors other than sole accessibility that impacts BMI, such as individual behavior and habits.

McFerran and Mukhopadhyay conducted a study on this behavioral aspect for predicting BMI. In their study, they studied people's personal beliefs on obesity and its cause. They assert that "such naive beliefs are important because they guide actual goal-directed behaviors" [3].

They found that people mainly believed that obesity is caused by either lack of exercise or by poor diet. Those holding the belief of the cause being lack of exercise were more likely to be overweight⁵ than those who held the belief that the cause is poor diet [3]. They explained that those who held the low-exercise belief tend to consume more food than those who held the poor-diet belief.

From these two studies, we see that the factors contributing to BMI are complex and varied, leaning towards more behavioral factors than concrete factors. Nonetheless, the density of various food sources may still influence food-acquiring habits and affect BMI.

3. Platforms Used

The project was conducted using mainly two cloud-based platforms. The data cleaning, transformation, and the training/testing of the model were performed on the Microsoft Azure Machine Learning Studio, and some data were also exported into Elasticsearch's cloud service for visualization.

3.1 Microsoft Azure Machine Studio

Since Azure ML Studio is a managed cloud platform, the exact specifications of the underlying hardware are not made public. However, depending on the service level, the service provides different levels of service from its service. This project was done on the free tier.

⁵ BMI > 24.9 is categorized as Overweight; >30 is Obese.

In the free tier, the service is limited compared to the paid tiers. Storage space, used for saving the datasets used in the project, is limited to 10GB. Each experiment has a maximum of 100 models, experiments are executed 1 node at a time, and experiments can run for no more than 1 hour. This means for larger and more complex projects, the execution may not complete within the time limit.

3.2 Elasticsearch Service

Like Azure, the Elasticsearch service is also a managed cloud service. Unlike Azure, a platform may be provisioned based on the desired configuration. The billing is dependant on the configuration of the desired instance that is to be provisioned. The configuration used for this project is billed at the rate of approximately \$0.02 cents an hour, with a monthly cost of around \$16 dollars.

The configuration used is 'gcp.data.highio.1', whose underlying machine specification is 12 cores, 78GB memory, and 3000GB storage. The instance used for the project is configured to use 1GB of memory and 30GB of storage on the machine.

4. Data Preparation

The data used is a 181.3MB set of 'csv' files obtained from the official USDA Economic Research Service website [1]. There are several sets of data within but there are only a few that contain the data relevant to the project.

4.1 Food Access

For obtaining features relating to food access, 'faps_access_puf.csv' was used. The data in this file is keyed by household number, 'hhnum', and contains a number of columns describing the number of food sources within a certain distance of the household of a specific type.

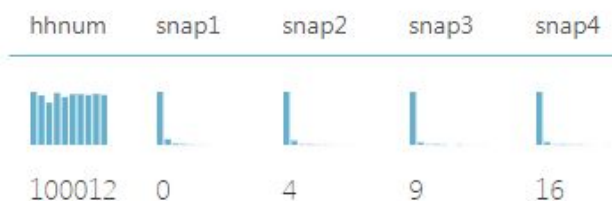


Figure 1. The first 5 columns of 'faps_access_puf.csv'

Figure 2. Breakdown of column names.

SNAP	SNAP-authorized retailers
SS	Superstores
SM	Supermarkets
CO	Combination Grocery
CS	Convenience Store
MLG	Medium or Large Grocery
FF	Fast Food Restaurants
NONFF	Non-fast Food Restaurants

Figure 3. Number component of column names.

1	within 0.25 miles
2	within 0.50 miles
3	within 1.00 miles
4	within 2.00 miles
5	within 5.00 miles
6	within 10.0 miles
7	within 15.0 miles
8	within 30.0 miles

For each column, the name denotes the type of food source, and the number denotes the corresponding distance away from the household within where the corresponding food sources are counted (e.g. 'FF5' meaning number of fast food restaurants within 5 miles). Each column name has 8 expanding radii, for a total of 64 features relating to food accessibility.

4.2 Individual Data

For obtaining the features and the future BMI label, the 'faps_individual_puf.csv' file was used. It is compositely keyed by household number ('hhnum') and the person number ('pnum') indicating the person within the household.

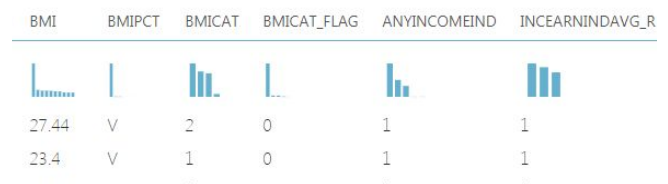


Figure 5. Example columns of 'faps_individual_puf.csv'

The columns contained describe various individual attributes of the associated survey participant. The column 'BMI' will primarily be used as the label for the model prediction.

In addition, there is a 'bmicat' column in the data with three values that denote the overweight status category: '1' denoting 'not overweight', '2' denoting 'overweight', and '3' denoting 'obese'. As a secondary objective, we will also try to classify this categorical column in a binary classification model — the column had been transformed with a Python script to change the 'bmicat' label to a binary class, '1' for not overweight, and '2' for overweight or obese.

4.3 Joining Food Access and Individual Data

Since both sets of data possess an 'hhnum' column, the two tables can be joined to produce a single table containing a row for every unique individual. Each row about the individual will possess the food accessibility features based on where they live — their household. The joining operation is accomplished using the Join Data module in Azure ML.

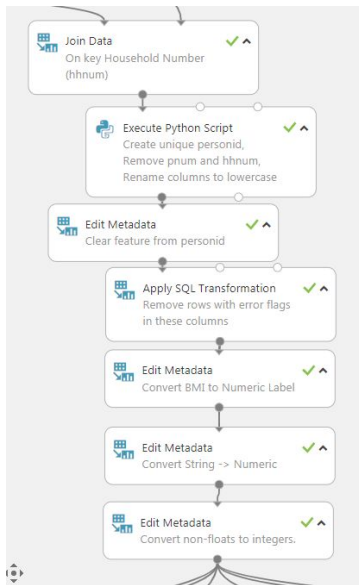


Figure 6. Series of data transformations.

After joining, a Python script was used to manipulate the tables to combine the 'hhnum' and 'pnun' fields into a single unique key by appending the 'pnun' to the 'hhnum' column⁶. The rest of the data cleaning was then performed, involving removing rows with errors and changing metadata, and the cleaned data was sent to the subsequent model training operations. Each model using this data splits the data into 66% training and 33% testing.

5. Results

In this project, several different configurations of prediction models were used. The various food accessibility columns will be the features used to predict the BMI label in a regression model. In the classification model, 'bmicat' will be used as the label, and the 'bmi' column will be removed to prevent it from influencing model performance since it should be reflective solely on the food accessibility factors⁷.

5.1 Regression Models

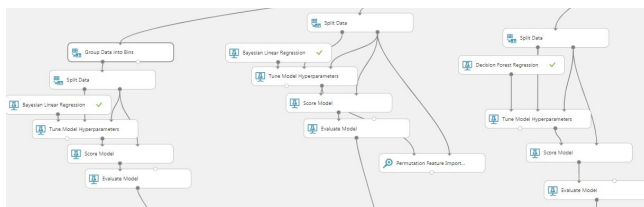


Figure 7. Binned and un-binned Bayesian Regression, and un-binned Decision Forest Regression.

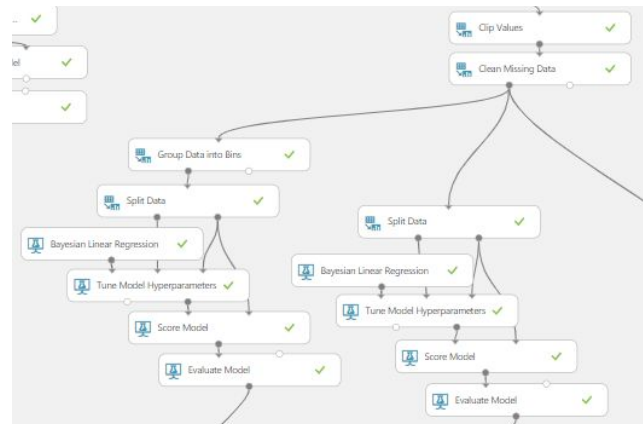


Figure 8. Binned and un-binned using values between 5th and 95th percentile.

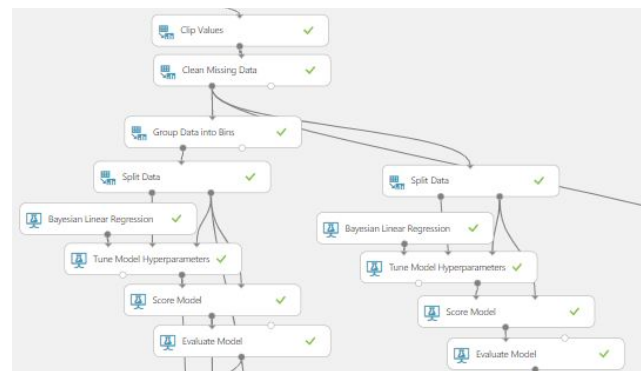


Figure 9. Binned and un-binned using values between 10th and 90th percentile.

Figure 10. Evaluation of regression models.

Type	Configuration	MAE	COD
Bayesian	Bin	3.3098	0.5605
Bayesian	No bin	3.3157	0.5607
Decision	No bin	3.2524	0.5734
Bayesian	Bin, 5-95%	2.3834	0.6061
Bayesian	No bin, 5-95%	2.3796	0.6061
Bayesian	Bin, 10-90%	2.0396	0.5959
Bayesian	No bin, 10-90%	2.0444	0.5951

The performance of the regression models is not very good as a whole. The low coefficient of determination across the board indicates that the label BMI is not very accurately predicted using the current labels. However, with average BMI being 25, a mean absolute error of approximately 2.0 is fairly accurate with the average error range being 7-10%.

5.2 Classification Models

For the classification models, the process is identical except that the model used is a Support Vector Machine model, and instead of the binning of features, cross-validation or lack thereof is used⁸. These models attempt to classify an

⁶ Scripts used available at the project Github repository: <https://github.com/shozonu/usda-foodaps-bmi-model>

⁷ Using a Python script, also available in the repository.

⁸ The modules are set up the same basic way in Azure ML as the regression models, and are visually similar.

The addition of frequent 'work', 'family' and 'store' keyword seem to also suggest that food-out events may coincide or may be in conjunction with other activities, such as store shopping with the family, or during work.

