

# Stat108\_FinalProject

## An Analysis on McDonalds' Menu and Their Nutrition Facts

Dylan Scoble and Aisha Lakshman

2/4/2022

### Introduction

In the 2004 documentary *Super Size Me*, writer and director Morgan Spurlock took on a month-long challenge to only eat McDonalds food. Spurlock experienced a multitude of health issues, including weight gain, cholesterol spike, and negative impacts on his energy and mood, demonstrating the fast-food chains' instrumental role in America's obesity epidemic (Stossel 2006). Spurlock's film not only emphasized the consequences of caloric intake, but also brought light to the nutritional attributes of McDonalds menu items that caused adverse health effects. There are many factors that impact the quality and quantity of calories, such as levels of fat, protein, and carbohydrates, which is why many dietitians support the notion that "not all calories are created equal" (Tolar-Peterson, 2021). Spurlock's documentary and existing literature inspired an investigation of McDonalds menu items' caloric and nutritional records.

Our research will address the following question: What nutritional attribute is most closely associated with calories for the McDonalds menu items? We will analyze a 2018 dataset from Kaggle titled "Nutritional Facts for McDonald's Menu" to answer our research question. Our chosen dataset provides nutritional information for all of McDonald's menu items, including calories, saturated fat, and cholesterol levels. Our research aims to guide the inspection of a nutritional label and to provide adequate information on the nutritional attributes that best estimates calories. Therefore, we will employ a modeling approach which estimates the closest association between calories and nutritional attributes.

To see which nutritional attribute is the best estimator for calories, we will create a linear model of each nutritional attribute with calories as the response variable. A linear model for regression analysis is useful in answering our question because it will allow us to confidently determine what nutritional attributes hold the closest association to calories.

### Our Data

The menu items and nutritional facts in our dataset are extracted from the McDonald's website (Kaggle 2016). The McDonalds Nutrition Calculator page provides information on how nutritional data were collected. According to McDonalds, existing nutritional data are "derived from testing conducted in accredited laboratories, published resources, or from information provided from McDonald's suppliers" (McDonalds 2017). The corporation also states that "% Daily Value" nutritional data is based on a 2,000 calorie diet (McDonalds 2017). The McDonald's nutritional calculator page provides the following message in fine print:

"All nutrition information is based on average values for ingredients and is rounded in accordance with current U.S. FDA NLEA regulations. Variation in serving sizes, preparation techniques, product testing and sources of supply, as well as regional and seasonal differences may affect the nutrition values for each product. In addition, product formulations change periodically. You should expect some variation in the nutrient content of the products purchased in our restaurants" (McDonalds 2017).

A sample of data can be viewed below.

[illegible]

We will also be removing all predictors that have “as % of Daily Value” attached at the end, since our purpose is not focused the daily values of the nutrients. These predictors add no value to our dataset or models.

In order to accurately compare the estimators, we want to standardize them first. To do this, we make sure every estimator is centered around a mean of zero with a variance of one. We create a new variable for each existing variable; this new variable contains the calculation of the formula  $\$ \text{new} = (\text{existing} - \text{mean}(\text{existing}) / \text{std\_dev}(\text{existing}))\$$

A sample of the final version of our datasets (food and beverage) can be viewed below.

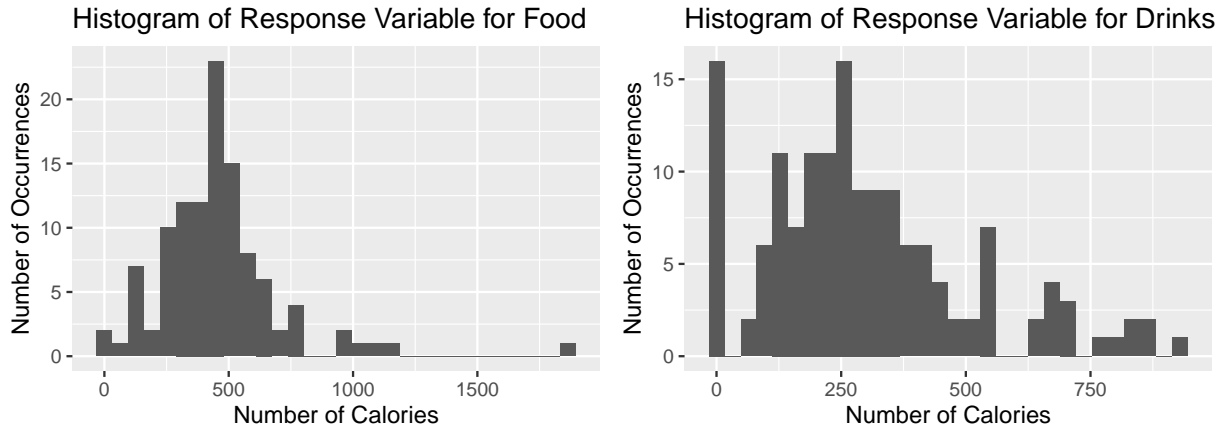
Food Data:

[illegible]

Beverage Data:

[illegible]

The next step is to create histograms for occurrences of food items (food\_data) and occurrences of drink items (bev\_data) against our response variable (calories). Based on the plots below, it seems that for both datasets, the Calories variable follows a normal distribution. A normal distribution tells us that data is centered around the mean and there are no obvious outliers. Additionally, a normal distribution is a good indicator of data independence.



The summary statistics for calories (mean, median, standard deviation, IQR) for food items and drink items is also displayed below. Notice that for both datasets, the mean is larger than the median. This is likely attributed to the fact that there are many items that McDonald’s advertises as “Zero Calories”. The large group of zero calorie items acts as an outlier for our normal distribution and brings the median down without affecting the mean as much.

Food Data:

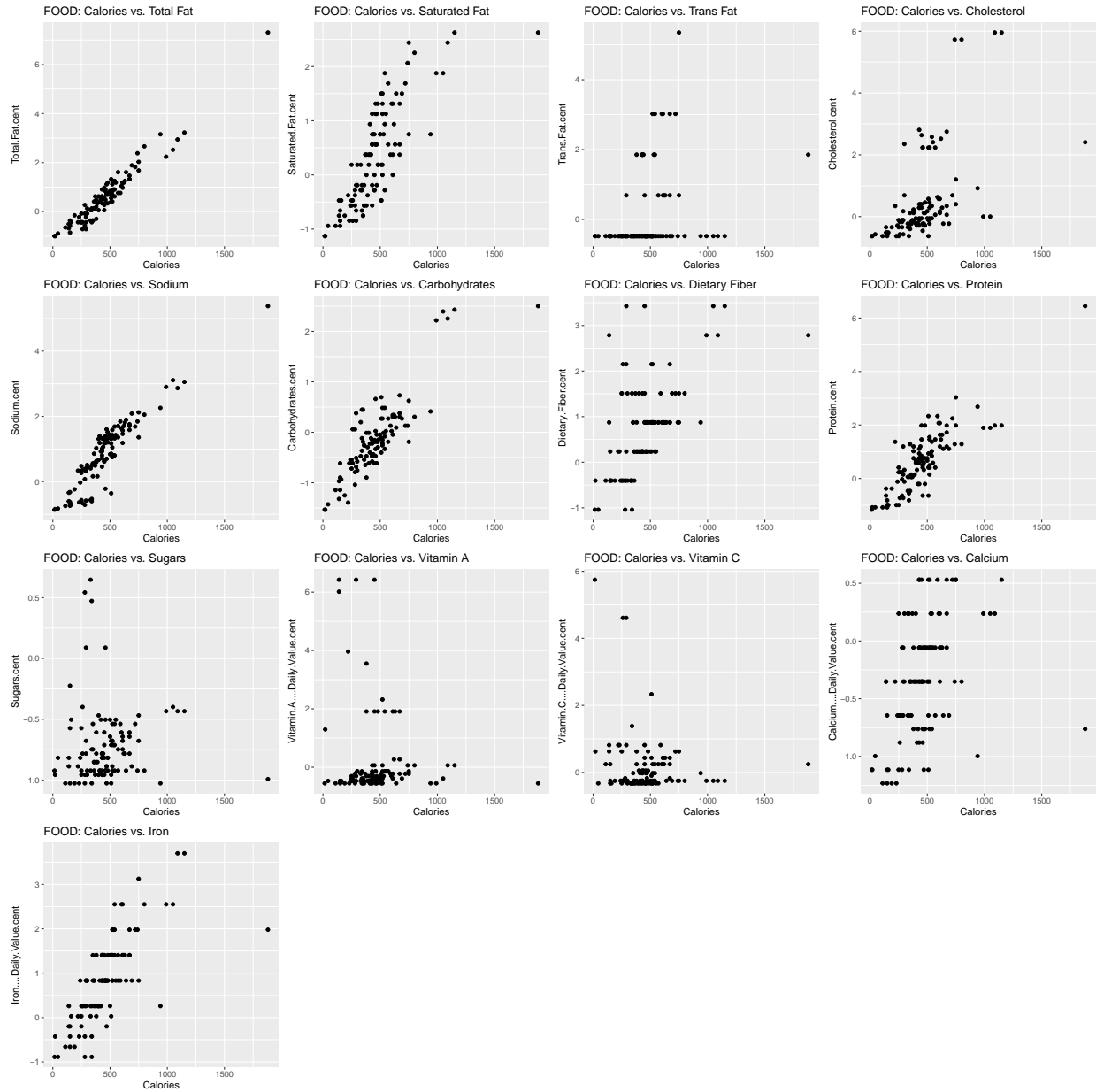
mean	median	std_dev	iqr
462.0909	445	249.3343	210

Beverage Data:

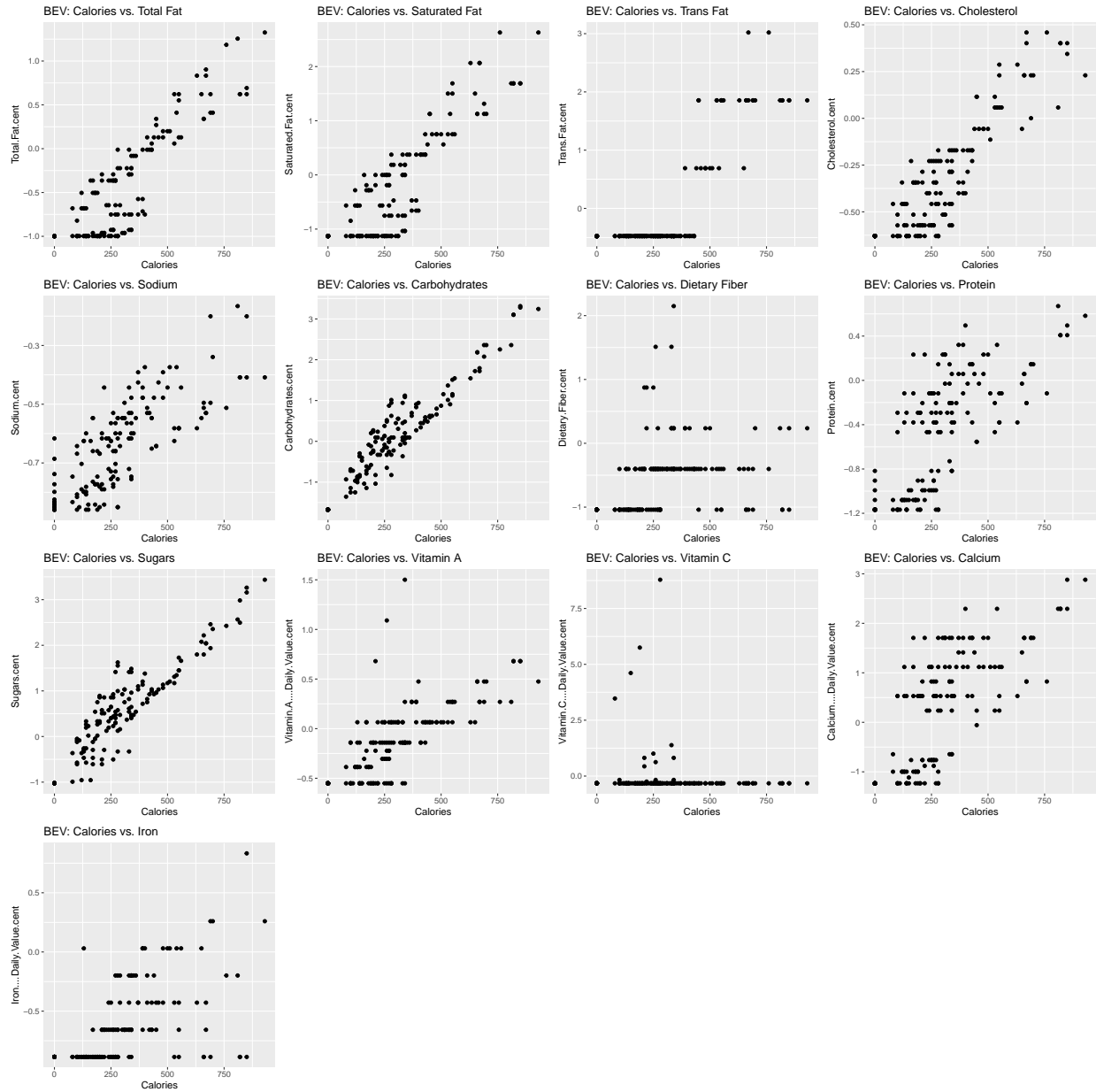
mean	median	std_dev	iqr
299.4667	270	208.8215	235

## Estimator Variables

We also want to view a plot of each variable’s relationship with our response. Below is every plot that describes the food dataset. We are creating a plot for each estimator (nutritional attribute) against our response variable (calories). These plots allow us to visualize the relationship between each estimator and calories for food menu items.



Now we will repeat this process with beverage menu items. Below is every plot that describes the beverage dataset. We are creating a plot for each estimator (nutritional attribute) against our response variable (calories). These plots allow us to visualize the relationship between each estimator and calories for beverage menu items.



## Model Selection with AIC

In order to conduct our experiment, we create a model that uses Calories as a response variable, and all centered variables as estimators.

The coefficients of our food model are displayed below.

term	estimate	std.error	statistic	p.value
(Intercept)	371.196	1.193	311.196	0.000
Total.Fat.cent	126.483	1.564	80.869	0.000
Saturated.Fat.cent	3.433	1.464	2.344	0.021
Trans.Fat.cent	0.563	0.656	0.858	0.393
Cholesterol.cent	-0.812	0.477	-1.703	0.092

term	estimate	std.error	statistic	p.value
Sodium.cent	-1.035	1.711	-0.605	0.547
Carbohydrates.cent	116.617	2.030	57.453	0.000
Dietary.Fiber.cent	-2.116	0.810	-2.614	0.010
Sugars.cent	-1.728	2.805	-0.616	0.539
Protein.cent	45.460	1.235	36.800	0.000
Vitamin.A...Daily.Value.cent	0.405	0.387	1.046	0.298
Vitamin.C...Daily.Value.cent	1.105	0.519	2.130	0.036
Calcium...Daily.Value.cent	-0.921	1.400	-0.657	0.512
Iron...Daily.Value.cent	-0.576	1.175	-0.490	0.625

In mathematical terms, we are creating an equation for a line where the x values are values of our estimators, and the y value is the response variable.

The equation for our food dataset can be read as the following:

$$\text{Calories} = 371.196 + 126.483(\text{Total.Fat}) + 3.433(\text{Saturated.Fat}) + 0.563(\text{Trans.Fat}) - 0.812(\text{Cholesterol}) - 1.035(\text{Sodium}) + 116.617(\text{Carbohydrates}) - 2.116(\text{Dietary.Fiber}) - 1.728(\text{Sugars}) + 45.460(\text{Protein}) + 0.405(\text{Vitamin.A....Daily.Value}) + 1.105(\text{Vitamin.C....Daily.Value}) - 0.921(\text{Calcium....Daily.Value}) - 0.576(\text{Iron....Daily.Value})$$

The coefficients of our food model are displayed below.

term	estimate	std.error	statistic	p.value
(Intercept)	356.609	5.746	62.060	0.000
Total.Fat.cent	103.013	8.553	12.044	0.000
Saturated.Fat.cent	15.861	6.730	2.357	0.020
Trans.Fat.cent	0.495	1.074	0.461	0.646
Cholesterol.cent	-0.848	9.999	-0.085	0.933
Sodium.cent	-8.990	11.414	-0.788	0.432
Carbohydrates.cent	117.309	3.585	32.719	0.000
Dietary.Fiber.cent	2.998	1.097	2.734	0.007
Sugars.cent	-8.659	3.617	-2.394	0.018
Protein.cent	35.835	6.878	5.210	0.000
Vitamin.A...Daily.Value.cent	3.869	1.917	2.018	0.046
Vitamin.C...Daily.Value.cent	1.712	0.481	3.556	0.001
Calcium...Daily.Value.cent	6.326	2.771	2.283	0.024
Iron...Daily.Value.cent	3.436	2.202	1.560	0.121

The equation for our food dataset can be read as the following:

$$\text{Calories} = 356.609 + 103.013(\text{Total.Fat}) + 15.861(\text{Saturated.Fat}) + 0.495(\text{Trans.Fat}) - 0.848(\text{Cholesterol}) - 8.990(\text{Sodium}) + 117.309(\text{Carbohydrates}) - 2.998(\text{Dietary.Fiber}) - 8.659(\text{Sugars}) + 35.835(\text{Protein}) + 3.869(\text{Vitamin.A....Daily.Value}) + 1.712(\text{Vitamin.C....Daily.Value}) + 6.326(\text{Calcium....Daily.Value}) - 3.436(\text{Iron....Daily.Value})$$

Because all of the variables have been standardized, they all have a mean of zero and a variance of one. This means that the coefficients for each model can be very easily compared with one another. Higher coefficients represent a greater relationship between the variable and the response, and lower coefficients represent a lesser relationship between the variable and the response.

Undoubtedly, the estimators with the greatest relationship to calorie count in both foods and beverages are Total Fat, Carbohydrates, and Protein. As an aside, we did not hypothesize that calcium would have a significant effect on beverages. Thinking back, this result makes sense because milk-based beverages are both high in calcium and calories.

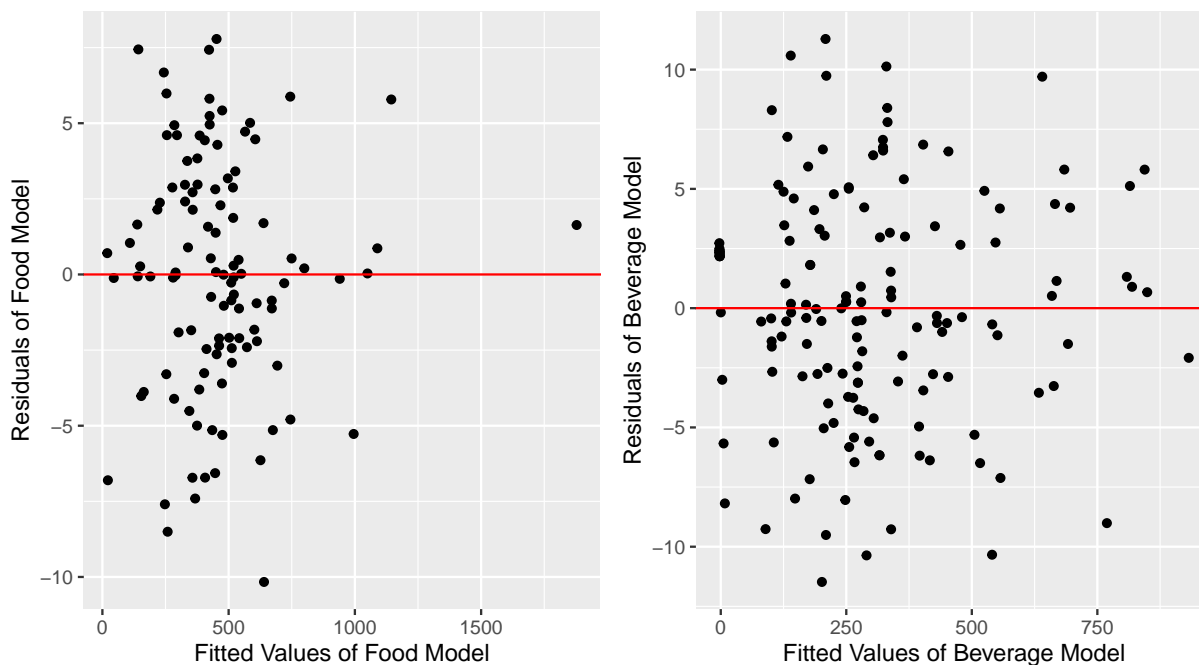
When first starting the analysis we decided to divide our data between food and drink items on the assumption that they have different caloric makeups. We split our dataset between food menu items (Breakfast, Beef & Pork, Chicken & Fish, Salads, Snacks & Sides, Desserts) and drink menu items (Coffee and Tea, Smoothies and Shakes, Beverages). For food menu items, we hypothesized total carbohydrates (grams) was the most accurate estimator of calories. For drink menu items, we hypothesized that total sugar (grams) was the most accurate estimator of calories. Upon completing our Exploratory Data Analysis and experiment, we concluded that best 3 estimators (nutritional attributes) were the same for food and beverage items. Thus, splitting the data did not lead to different results. If the same experiment were performed without splitting the dataset, we can expect Total Fat, Carbohydrates, and Protein to remain the highest coefficients.

## Checking Assumptions

Before accepting our results as truth, we must inspect where the data and the model come from. There are certain assumptions we were accepting as truth when conducting our experiment. In order to confirm our results, we must verify that our prior assumptions are correct.

### Constant Variance

To determine if the constant variance assumption is satisfied, we can inspect a plot of fitted values versus the residuals for food data and beverage data. The figures below plot fitted values on the x-axis and the residuals of fitted values along the y-axis. The constant variance assumption is satisfied if residual points are equally distributed with regard to the red line (residuals = 0) and if there is no clear shape or pattern in the fitted values versus residual plot. The plots below indicate that the constant variance assumption is satisfied for the food model and the beverage model.



«««< HEAD

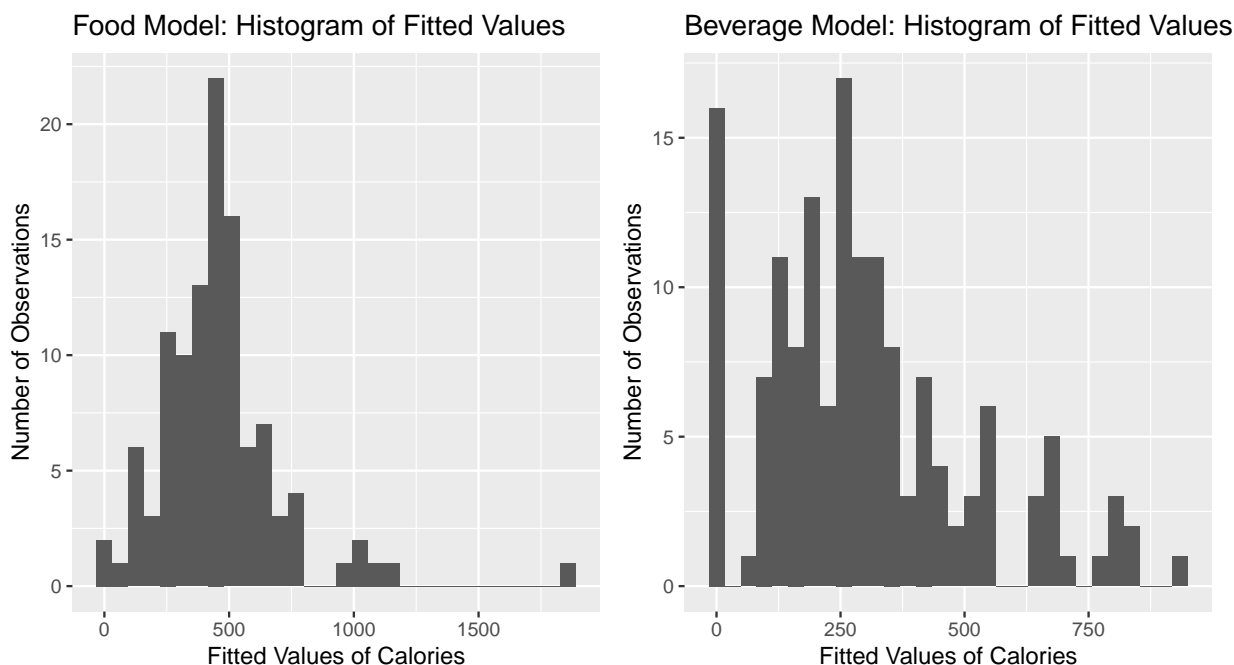
===== »»»> 354450a1af3a76cdf6ad85af4a6d045b64499fe9

### Linearity

The plots on page 4 and 5 show that for both food and beverages, there are linear relationships between calories and other nutritional attributes. This tells us that the linearity assumption is satisfied.

## Normality

To check if the normality assumption is satisfied or not, we can inspect the distribution of the fitted values for calories. To do so, we will create a histogram of fitted values against the number of observations for the food model and beverage model. The normality assumption is satisfied if the histogram follows a normal distribution. The plots below indicate that the normality distribution is satisfied for both of our models. The histogram for the food model clearly follows a normal distribution. The beverage data loosely follows a normal distribution, with the exception of a few outliers.



## Independence

A crucial assumption of linear regression is the independence of observations. Looking at how our data was collected will indicate if the independence assumption is satisfied or not. Given that our dataset consists of nutritional attributes for each McDonalds menu item, each observation is independent. A menu item's observed nutritional attributes does not rely on other menu items. In part by FDA Menu Labeling Requirements (2020) the process by which our data was collected ensures data validity and that we are working with a random sample.

## Conclusion

Throughout this analysis, we answered the question: What nutritional attribute is the best estimator for a McDonald's menu item's calories? In other words, we sought to find the nutritional attribute that has the closest relationship with our response variable, which is calories.

We used a dataset directly from the McDonald's corporation that contains the nutritional facts of every menu item. Then, we split our data into two datasets, food and beverages, and standardized each feature. We created a model that compares the features against each other. Undoubtedly, the three greatest estimators of calorie count in both food and beverages was total fat, carbohydrates, and protein. Our hypothesis regarding foods was partially correct; carbohydrates is a major predictor of a food's calories. However, we did not hypothesize that fat content and would have the impact it did, since carbohydrates have the reputation of being the most satiating. For beverages, our hypothesis was entirely wrong. Due to the nature of diet sodas, we hypothesized that a drink's sugar content would closely relate to its calorie count. We also thought that



foods and beverages would have vastly different results. We are grateful that this study was able to challenge our biases and stigmas regarding certain foods and food groups.

Although this study was done specifically on the McDonald's menu, we hope that the results of this experiment helps readers understand where their calories come from. On top of this, we hope to make reading a nutritional label easier for everyone. In the case where nutritional information is not accessible, consumers can still make educated inferences of a food's calorie count based on the results of this study.

## References

Tolar-Peterson, Terezia. 2021. "Not all calories are created equal - a dietician explains the different ways the kinds of foods you eat matter to your body". The Conversation. Retrieved February 8th, 2022. <https://theconversation.com/not-all-calories-are-equal-a-dietitian-explains-the-different-ways-the-kinds-of-foods-you-eat-matter-to-your-body-156900>

Stossel, John. 2006. "'Super Size Me' Carries Weight With Critics". ABC News. Retrieved February 8th, 2022. [https://docs.google.com/document/d/1XB-22QylvnbasBKe7n\\_DfkkENcZWRvIR5X8vZgOK6LY/edit#](https://docs.google.com/document/d/1XB-22QylvnbasBKe7n_DfkkENcZWRvIR5X8vZgOK6LY/edit#)

Food and Drug Administration. 2020. "Menu and Food Labeling Requirements". <https://www.fda.gov/food/food-labeling-nutrition/menu-labeling-requirements>

In text: McDonalds (2017) <https://www.mcdonalds.com/us/en-us/about-our-food/nutrition-calculator.html>

FDA (2014) <https://www.fda.gov/nutrition-labeling-and-education-act-nlea-requirements-attachment-1>