

Stat108_FinalProject

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 dieticians 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).

The Nutrition Labeling and Education Act (NLEA) requires corporations like McDonald's to provide adequate nutritional information and labeling for menu items (FDA 2014). FDA NLEA guidelines provide clear and consistent nutritional labeling requirements, allowing consumers to make informed dietary choices for themselves and their loved ones.

[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.

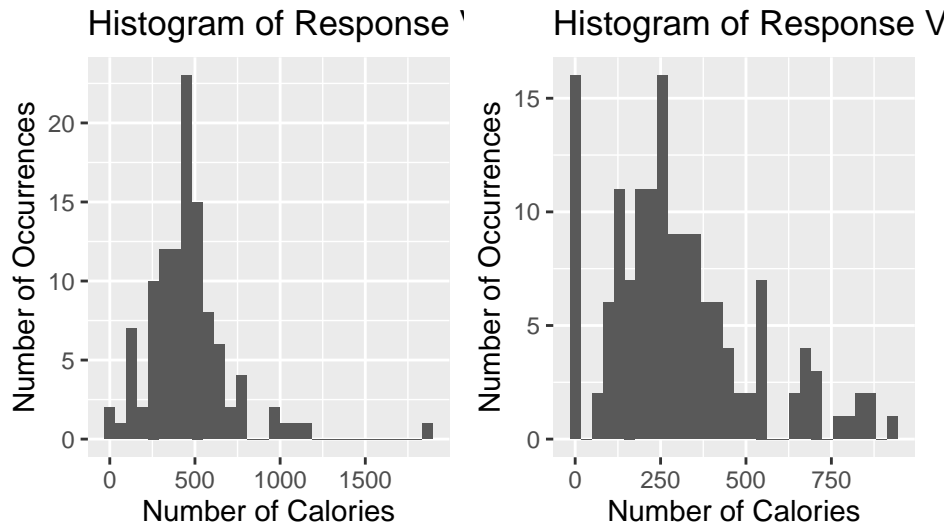
In order to accurately compare the estimators, we want to standardize them first. To do this, we make sure every estimator is centered around a mean of zero with a variance of one. We create a new variable for each existing variable; this new variable contains the calculation of the formula $\text{new}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

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

[illegible]

Variable	Unit	2000-2001		2001-2002		2002-2003		2003-2004		2004-2005		2005-2006		2006-2007		2007-2008		2008-2009		2009-2010		2010-2011		2011-2012		2012-2013		2013-2014		2014-2015		2015-2016		2016-2017		2017-2018		2018-2019		2019-2020		2020-2021		2021-2022		2022-2023		2023-2024		2024-2025		2025-2026		2026-2027		2027-2028		2028-2029		2029-2030		2030-2031		2031-2032		2032-2033		2033-2034		2034-2035		2035-2036		2036-2037		2037-2038		2038-2039		2039-2040		2040-2041		2041-2042		2042-2043		2043-2044		2044-2045		2045-2046		2046-2047		2047-2048		2048-2049		2049-2050		2050-2051		2051-2052		2052-2053		2053-2054		2054-2055		2055-2056		2056-2057		2057-2058		2058-2059		2059-2060		2060-2061		2061-2062		2062-2063		2063-2064		2064-2065		2065-2066		2066-2067		2067-2068		2068-2069		2069-2070		2070-2071		2071-2072		2072-2073		2073-2074		2074-2075		2075-2076		2076-2077		2077-2078		2078-2079		2079-2080		2080-2081		2081-2082		2082-2083		2083-2084		2084-2085		2085-2086		2086-2087		2087-2088		2088-2089		2089-2090		2090-2091		2091-2092		2092-2093		2093-2094		2094-2095		2095-2096		2096-2097		2097-2098		2098-2099		2099-2100		2100-2101		2101-2102		2102-2103		2103-2104		2104-2105		2105-2106		2106-2107		2107-2108		2108-2109		2109-2110		2110-2111		2111-2112		2112-2113		2113-2114		2114-2115		2115-2116		2116-2117		2117-2118		2118-2119		2119-2120		2120-2121		2121-2122		2122-2123		2123-2124		2124-2125		2125-2126		2126-2127		2127-2128		2128-2129		2129-2130		2130-2131		2131-2132		2132-2133		2133-2134		2134-2135		2135-2136		2136-2137		2137-2138		2138-2139		2139-2140		2140-2141		2141-2142		2142-2143		2143-2144		2144-2145		2145-2146		2146-2147		2147-2148		2148-2149		2149-2150		2150-2151		2151-2152		2152-2153		2153-2154		2154-2155		2155-2156		2156-2157		2157-2158		2158-2159		2159-2160		2160-2161		2161-2162		2162-2163		2163-2164		2164-2165		2165-2166		2166-2167		2167-2168		2168-2169		2169-2170		2170-2171		2171-2172		2172-2173		2173-2174		2174-2175		2175-2176		2176-2177		2177-2178		2178-2179		2179-2180		2180-2181		2181-2182		2182-2183		2183-2184		2184-2185		2185-2186		2186-2187		2187-2188		2188-2189		2189-2190		2190-2191		2191-2192		2192-2193		2193-2194		2194-2195		2195-2196		2196-2197		2197-2198		2198-2199		2199-2200		2200-2201		2201-2202		2202-2203		2203-2204		2204-2205		2205-2206		2206-2207		2207-2208		2208-2209		2209-2210		2210-2211		2211-2212		2212-2213		2213-2214		2214-2215		2215-2216		2216-2217		2217-2218		2218-2219		2219-2220		2220-2221		2221-2222		2222-2223		2223-2224		2224-2225		2225-2226	
----------	------	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--	-----------	--

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. ^ AISHA can u add a sentence why the normal distribution is a good thing



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

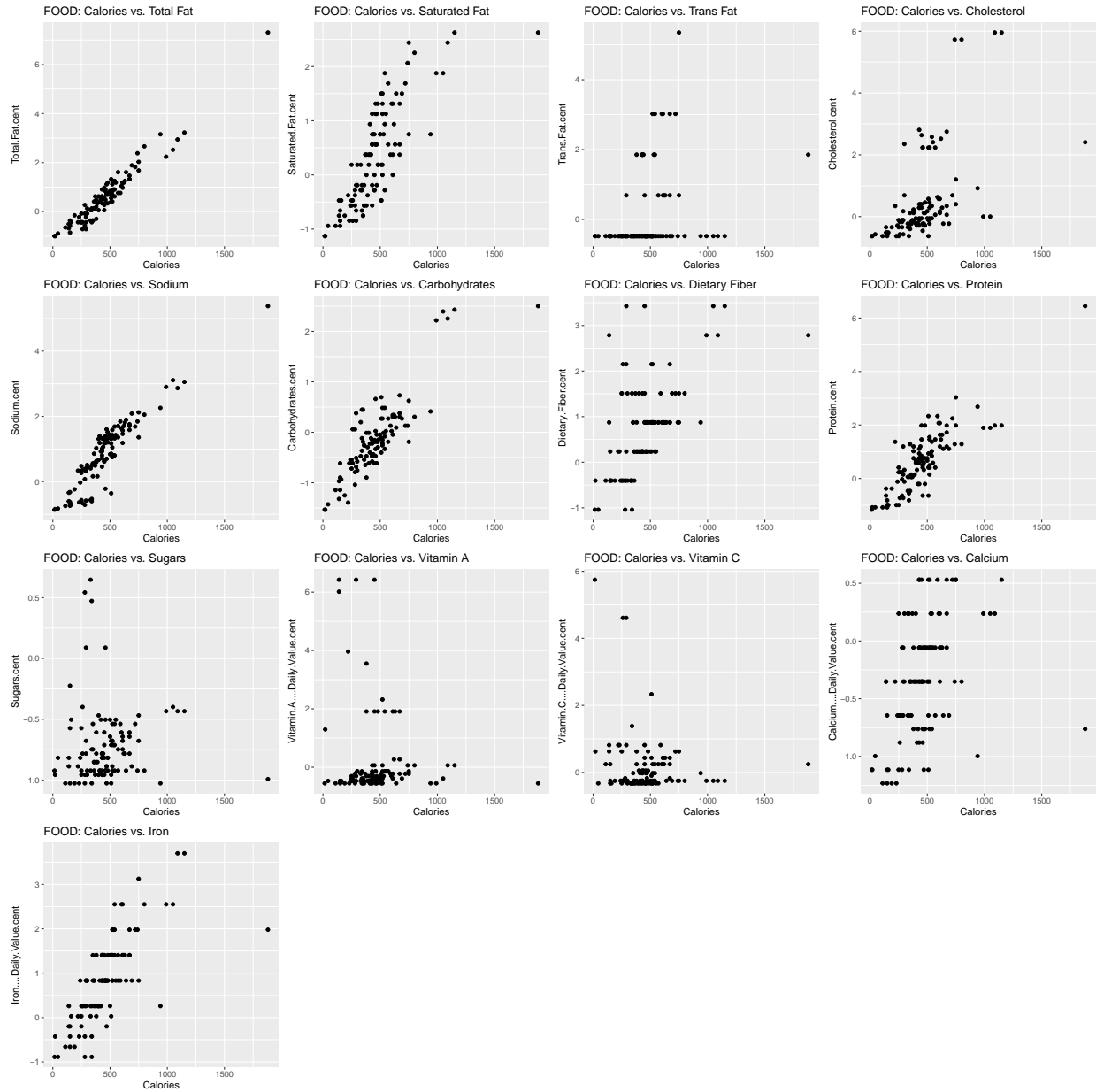
Beverage Data:

```
kable(t2, format = "markdown")
```

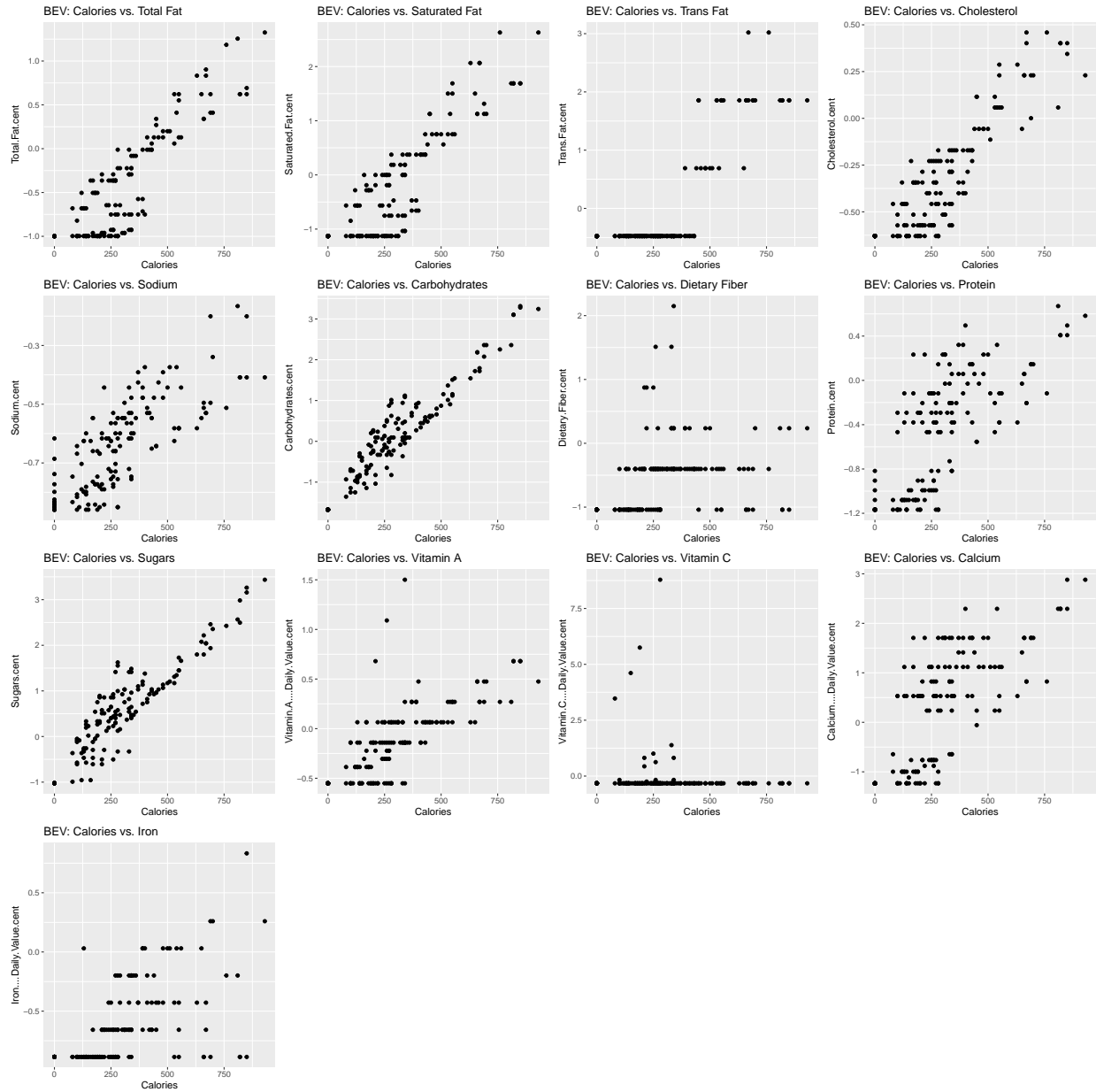
mean	median	std_dev	iqr
299.4667	270	208.8215	235

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. ^ AISHA can you add some here that talks about what the plots are showing and why we are showing them



Below is every plot that describes the beverage dataset. ^ AISHA can you add some here that talks about what the plots are showing and why we are showing them



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.

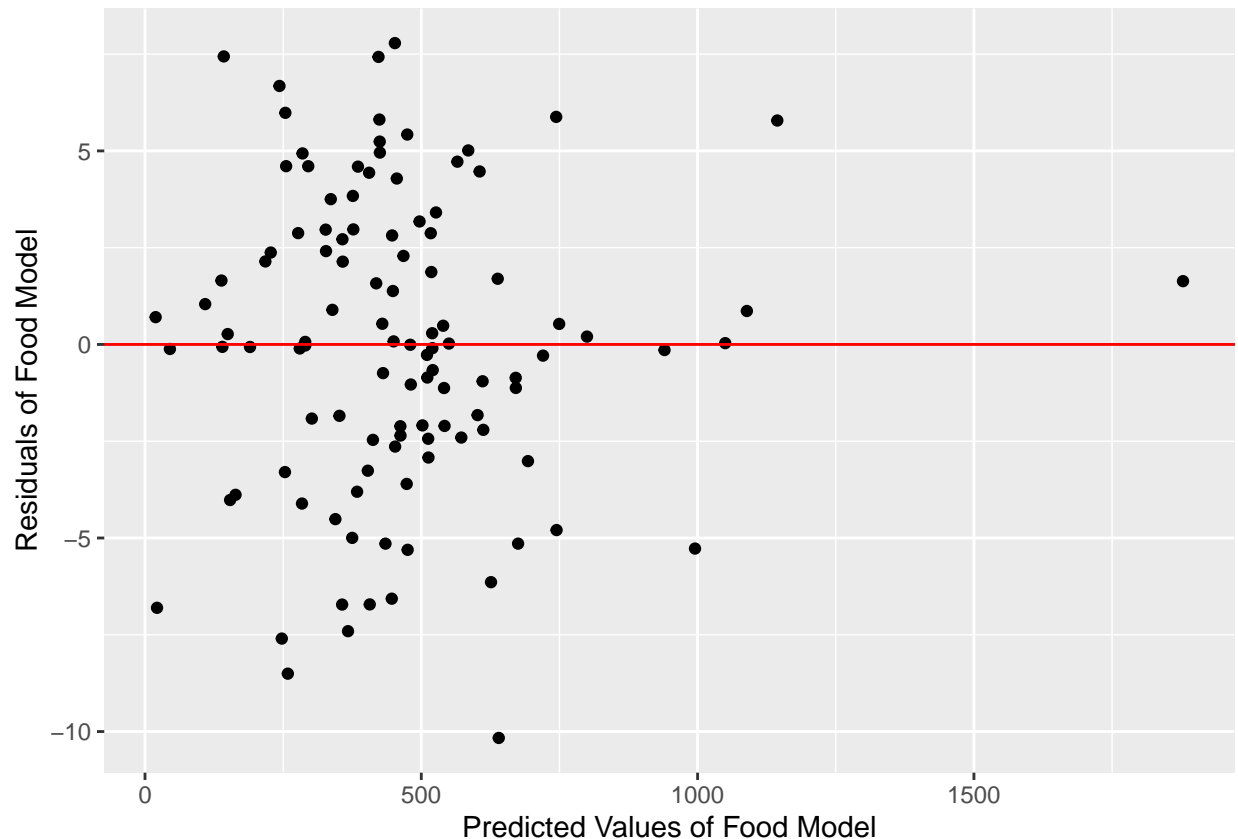
Undoubtedly, the estimators with the greatest relationship to calorie count in both foods and beverages are Total Fat, Carbohydrates, and Protein. As an aside, we did not hypothesize that calcium would have a significant effect on beverages. Thinking back, this result makes sense because milk-based beverages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

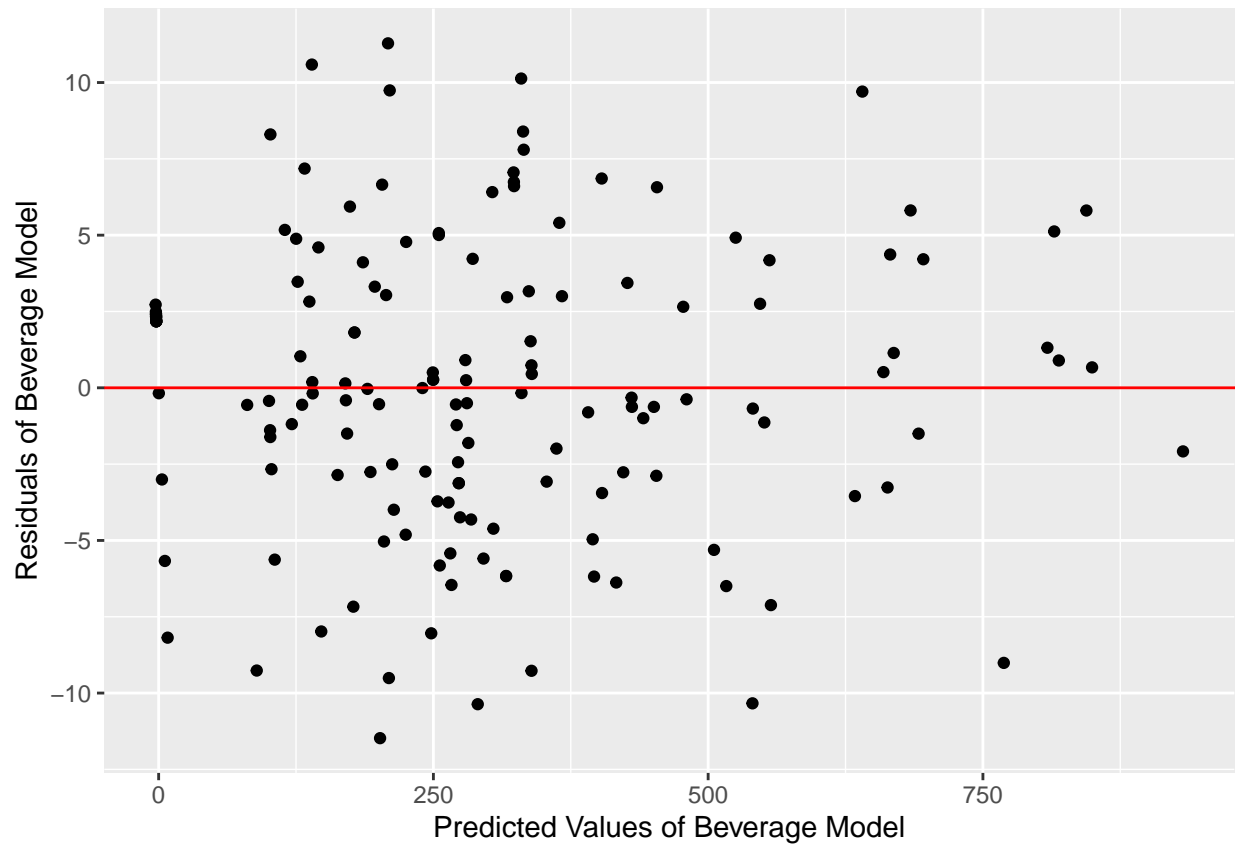
When first starting the analysis we decided to divide our data between food and drink items on the assumption that they have different caloric makeups. We split our dataset between food menu items (Breakfast, Beef & Pork, Chicken & Fish, Salads, Snacks & Sides, Desserts) and drink menu items (Coffee 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 AISHA can you write about how these plots were created and explain that the constant variance assumption is satisfied



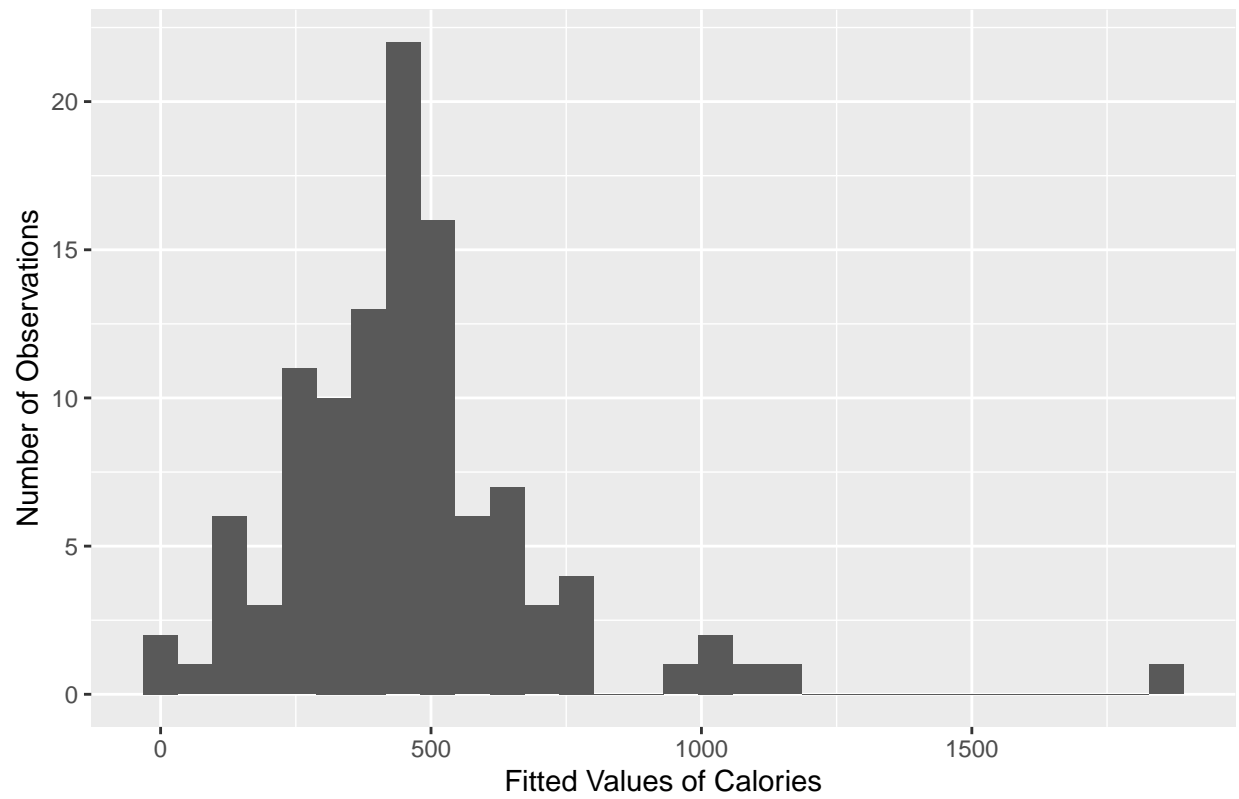


Linearity

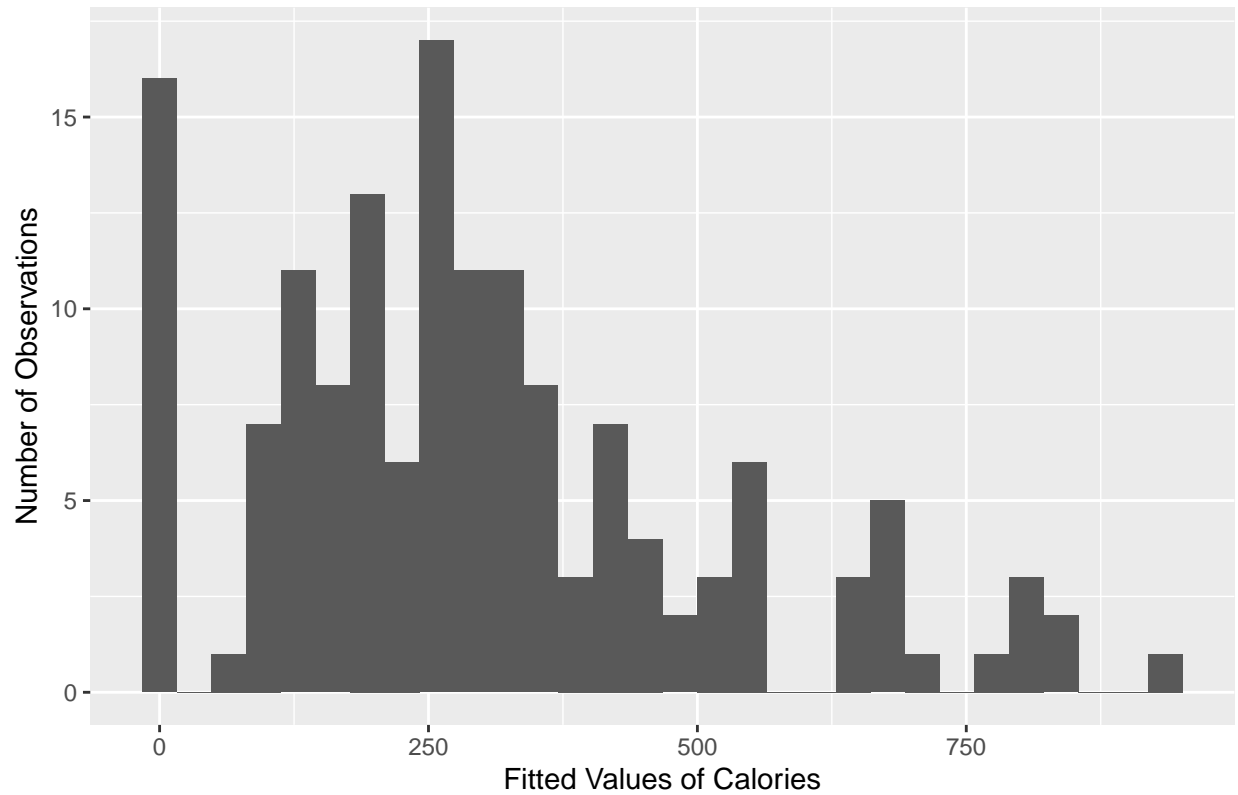
The plots on page 4 and 5 show that for both food and beverages, there are linear relationships between calories and other nutritional attributes. This tells us that the linearity assumption is satisfied.

Normality AISHA can you explain the same for the normality assumption

Food Model: Histogram of Fitted Values



Beverage Model: Histogram of Fitted Values



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.