

# Stat108\_FinalProject

Dylan Scoble and Aisha Lakshman

2/4/2022

«««< HEAD:Analysis.Rmd

## 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 the best predictor of calories for the McDonalds menu items? Since drink and food items have quite different caloric makeups, we are dividing 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 hypothesize total carbohydrates (grams) is the most accurate predictor of calories. For drink menu items, we hypothesize that total sugar (grams) is the most accurate predictor of calories. 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. We will create a linear model for each nutritional attribute with calories as the response variable for each predictor. A linear model for regression analysis is useful in answering our question because it will allow us to confidently determine what nutritional attributes matter the most for calories and predict an item's calorie count based on its predictors.

## References

- Tolar-Peterson, Terezie. 2021. "Not all calories are created equal - a dietitian 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#)

## Our Data

```
data <- read.csv("data/menu 2.csv")
glimpse(data)

## Rows: 260
## Columns: 24
## $ Category      <chr> "Breakfast", "Breakfast", "Breakfast", "~
## $ Item           <chr> "Egg McMuffin", "Egg White Delight", "Sa~
## $ Serving.Size   <chr> "4.8 oz (136 g)", "4.8 oz (135 g)", "3.9~
## $ Calories       <int> 300, 250, 370, 450, 400, 430, 460, 520, ~
## $ Calories.from.Fat <int> 120, 70, 200, 250, 210, 210, 230, 270, 1~
## $ Total.Fat      <dbl> 13, 8, 23, 28, 23, 23, 26, 30, 20, 25, 2~
## $ Total.Fat....Daily.Value. <int> 20, 12, 35, 43, 35, 36, 40, 47, 32, 38, ~
## $ Saturated.Fat  <dbl> 5, 3, 8, 10, 8, 9, 13, 14, 11, 12, 12, 1~
## $ Saturated.Fat....Daily.Value. <int> 25, 15, 42, 52, 42, 46, 65, 68, 56, 59, ~
## $ Trans.Fat      <dbl> 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, ~
## $ Cholesterol    <int> 260, 25, 45, 285, 50, 300, 250, 250, 35,~
## $ Cholesterol....Daily.Value. <int> 87, 8, 15, 95, 16, 100, 83, 83, 11, 11, ~
## $ Sodium         <int> 750, 770, 780, 860, 880, 960, 1300, 1410~
## $ Sodium....Daily.Value. <int> 31, 32, 33, 36, 37, 40, 54, 59, 54, 59, ~
## $ Carbohydrates  <int> 31, 30, 29, 30, 30, 31, 38, 43, 36, 42, ~
## $ Carbohydrates....Daily.Value. <int> 10, 10, 10, 10, 10, 10, 13, 14, 12, 14, ~
## $ Dietary.Fiber  <int> 4, 4, 4, 4, 4, 4, 2, 3, 2, 3, 2, 3, ~
## $ Dietary.Fiber....Daily.Value. <int> 17, 17, 17, 17, 17, 18, 7, 12, 7, 12, 6,~
## $ Sugars         <int> 3, 3, 2, 2, 2, 3, 3, 4, 3, 4, 2, 3, 2, ~
## $ Protein        <int> 17, 18, 14, 21, 21, 26, 19, 19, 20, 20, ~
## $ Vitamin.A....Daily.Value. <int> 10, 6, 8, 15, 6, 15, 10, 15, 2, 6, 0, 4,~
## $ Vitamin.C....Daily.Value. <int> 0, 0, 0, 0, 0, 2, 8, 8, 8, 8, 0, 0, 0, ~
## $ Calcium....Daily.Value. <int> 25, 25, 25, 30, 25, 30, 15, 20, 15, 15, ~
## $ Iron....Daily.Value. <int> 15, 8, 10, 15, 10, 20, 15, 20, 10, 15, 1~
```

## Exploratory Data Analysis

For this project, we understand that foods and beverages may have different predictors for their number of calories. Therefore, we will be splitting our dataset into two different dataframes: one for foods, and one for beverages.

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.

```
data <- data %>%
  select(Category, Item, Serving.Size, Calories, Calories.from.Fat, Total.Fat,
         Saturated.Fat, Trans.Fat, Cholesterol, Sodium, Carbohydrates,
         Dietary.Fiber, Sugars, Protein, Vitamin.A....Daily.Value.,
         Vitamin.C....Daily.Value., Calcium....Daily.Value., Iron....Daily.Value.)
glimpse(data)
```

```
## Rows: 260
```

```
## Columns: 18
## $ Category      <chr> "Breakfast", "Breakfast", "Breakfast", "Brea~
## $ Item           <chr> "Egg McMuffin", "Egg White Delight", "Sausag~
## $ Serving.Size   <chr> "4.8 oz (136 g)", "4.8 oz (135 g)", "3.9 oz ~
## $ Calories       <int> 300, 250, 370, 450, 400, 430, 460, 520, 410, ~
## $ Calories.from.Fat <int> 120, 70, 200, 250, 210, 210, 230, 270, 180, ~
## $ Total.Fat      <dbl> 13, 8, 23, 28, 23, 23, 26, 30, 20, 25, 27, 3~
## $ Saturated.Fat  <dbl> 5, 3, 8, 10, 8, 9, 13, 14, 11, 12, 12, 13, 1~
## $ Trans.Fat      <dbl> 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, ~
## $ Cholesterol    <int> 260, 25, 45, 285, 50, 300, 250, 250, 35, 35, ~
## $ Sodium         <int> 750, 770, 780, 860, 880, 960, 1300, 1410, 13~
## $ Carbohydrates  <int> 31, 30, 29, 30, 30, 31, 38, 43, 36, 42, 34, ~
## $ Dietary.Fiber  <int> 4, 4, 4, 4, 4, 4, 2, 3, 2, 3, 2, 3, 2, ~
## $ Sugars         <int> 3, 3, 2, 2, 2, 3, 3, 4, 3, 4, 2, 3, 2, 3, ~
## $ Protein        <int> 17, 18, 14, 21, 21, 26, 19, 19, 20, 20, 11, ~
## $ Vitamin.A....Daily.Value. <int> 10, 6, 8, 15, 6, 15, 10, 15, 2, 6, 0, 4, 6, ~
## $ Vitamin.C....Daily.Value. <int> 0, 0, 0, 0, 0, 2, 8, 8, 8, 8, 0, 0, 0, 0, ~
## $ Calcium....Daily.Value.  <int> 25, 25, 25, 30, 25, 30, 15, 20, 15, 15, 6, 8~
## $ Iron....Daily.Value.    <int> 15, 8, 10, 15, 10, 20, 15, 20, 10, 15, 15, 1~
```

```
count(data, Category)
```

```
##           Category    n
## 1      Beef & Pork   15
## 2      Beverages   27
## 3      Breakfast   42
## 4  Chicken & Fish   27
## 5  Coffee & Tea    95
## 6      Desserts     7
## 7      Salads       6
## 8 Smoothies & Shakes 28
## 9  Snacks & Sides   13
```

```
food_data <- data %>%
  filter(Category == "Beef & Pork" |
         Category == "Breakfast" |
         Category == "Chicken & Fish" |
         Category == "Desserts" |
         Category == "Salads" |
         Category == "Snacks & Sides")

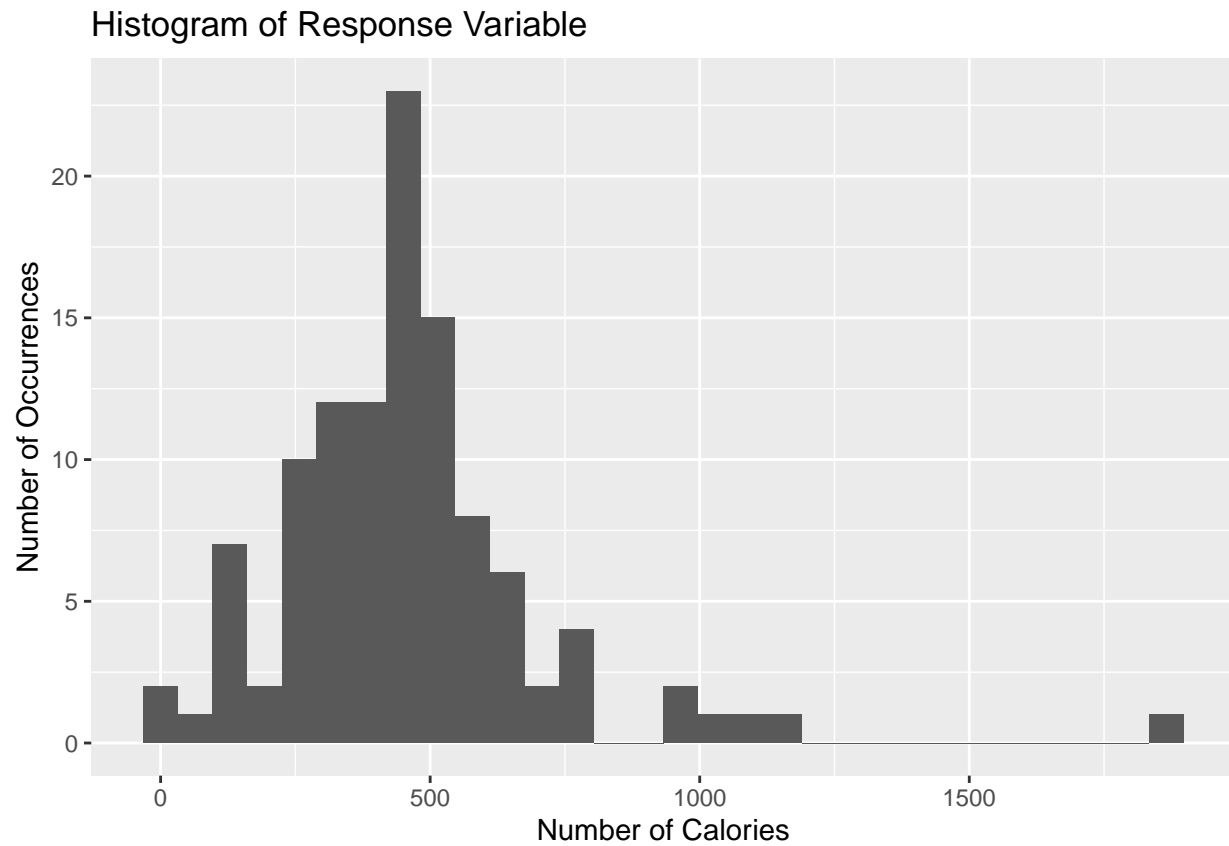
bev_data <- data %>%
  filter(Category == "Beverages" |
         Category == "Coffee & Tea" |
         Category == "Smoothies & Shakes")
```

## Response Variable

```
ggplot(data=food_data, aes(x=Calories)) +
  geom_histogram() +
  labs(title="Histogram of Response Variable",
```

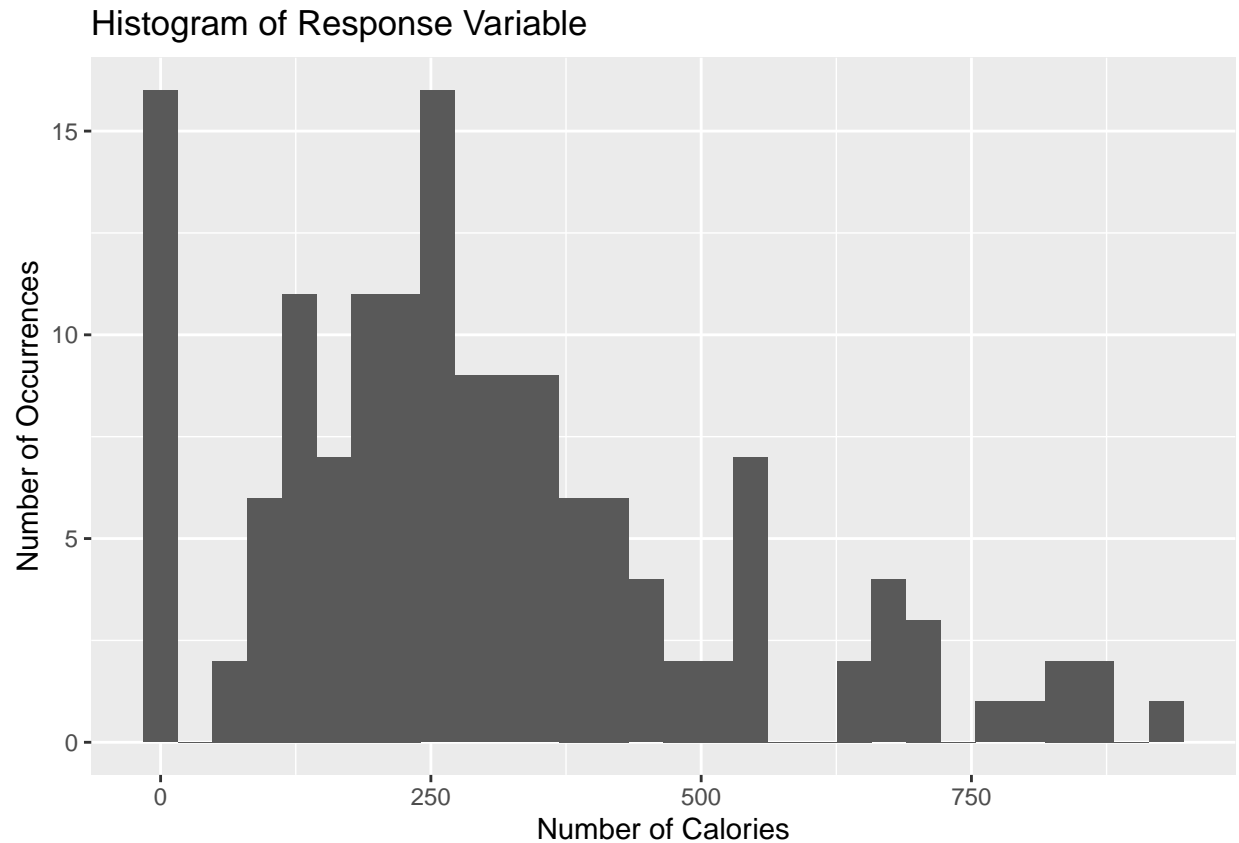
```
x="Number of Calories",  
y="Number of Occurrences")
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
ggplot(data=bev_data, aes(x=Calories)) +  
  geom_histogram() +  
  labs(title="Histogram of Response Variable",  
        x="Number of Calories",  
        y="Number of Occurrences")
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
food_data %>%
  summarise(mean = mean(Calories),
            median = median(Calories),
            std_dev = sd(Calories),
            iqr = IQR(Calories))
```

```
##      mean median  std_dev iqr
## 1 462.0909   445 249.3343 210
```

```
bev_data %>%
  summarise(mean = mean(Calories),
            median = median(Calories),
            std_dev = sd(Calories),
            iqr = IQR(Calories))
```

```
##      mean median  std_dev iqr
## 1 299.4667   270 208.8215 235
```

## Regression

```
food_model = lm(Calories ~ Total.Fat + Saturated.Fat + Trans.Fat + Cholesterol + Sodium + Carbohydrates
  Dietary.Fiber + Sugars + Protein + Vitamin.A....Daily.Value. +
```

```
Vitamin.C....Daily.Value. + Calcium....Daily.Value. + Iron....Daily.Value., data = food_data)

tidy(food_model, conf.int = TRUE) %>%
  kable(format="markdown", digits = 5)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-1.12936	1.19051	-0.94864	0.34519	-3.49250	1.23377
Total.Fat	8.90349	0.11010	80.86937	0.00000	8.68495	9.12203
Saturated.Fat	0.64499	0.27515	2.34419	0.02113	0.09883	1.19115
Trans.Fat	1.31130	1.52854	0.85788	0.39310	-1.72284	4.34544
Cholesterol	-0.00930	0.00546	-1.70272	0.09186	-0.02015	0.00154
Sodium	-0.00179	0.00297	-0.60505	0.54657	-0.00768	0.00409
Carbohydrates	4.12772	0.07185	57.45281	0.00000	3.98511	4.27033
Dietary.Fiber	-1.35003	0.51653	-2.61367	0.01040	-2.37533	-0.32473
Sugars	-0.06024	0.09780	-0.61593	0.53940	-0.25436	0.13389
Protein	3.97862	0.10811	36.80011	0.00000	3.76402	4.19323
Vitamin.A...Daily.Value.	0.01661	0.01587	1.04647	0.29797	-0.01489	0.04811
Vitamin.C...Daily.Value.	0.04194	0.01969	2.13042	0.03569	0.00286	0.08102
Calcium...Daily.Value.	-0.05409	0.08228	-0.65740	0.51250	-0.21743	0.10924
Iron...Daily.Value.	-0.06604	0.13469	-0.49029	0.62505	-0.33339	0.20132

```
bev_model = lm(Calories ~ Total.Fat + Sodium + Carbohydrates + Sugars + Protein + Vitamin.A....Daily.Value.
  Vitamin.C....Daily.Value. + Calcium....Daily.Value. + Iron....Daily.Value., data = bev_data)

tidy(bev_model, conf.int = TRUE) %>%
  kable(format="markdown", digits = 5)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	-1.07350	0.94628	-1.13444	0.25855	-2.94435	0.79735
Total.Fat	9.04943	0.08198	110.38674	0.00000	8.88735	9.21151
Sodium	-0.05167	0.01732	-2.98320	0.00337	-0.08591	-0.01743
Carbohydrates	4.34711	0.11325	38.38434	0.00000	4.12321	4.57102
Sugars	-0.47748	0.11734	-4.06924	0.00008	-0.70946	-0.24549
Protein	3.79427	0.56608	6.70269	0.00000	2.67510	4.91345
Vitamin.A...Daily.Value.	0.15757	0.07876	2.00069	0.04736	0.00186	0.31328
Vitamin.C...Daily.Value.	0.04432	0.01808	2.45206	0.01543	0.00859	0.08006
Calcium...Daily.Value.	0.26133	0.15820	1.65187	0.10080	-0.05145	0.57411
Iron...Daily.Value.	0.76536	0.21183	3.61317	0.00042	0.34657	1.18416