

Explore the Hot vs. Cold Cereals Data Set Using R and R Studio

Didem Bulut Aykurt

MIS510-1 – Data Mining and Visualization

Colorado State University-Global Campus

Dr. Emmanuel Tsukerman

Now 27,2022

Module 2-Option 1

This report is a data analysis of the Cereal dataset. I aim to write an R code for summary statistics such as mean, standard deviation, min, max, median, length, and the sum of missing values. Additionally, create a histogram and box plot.

The name of the dataset is Cereals dataset describes hot and cold cereals containing each cereal's nutrition and 77 different cereal products—the data sources from CSU-Global. The dataset has 77 observations and 16 variables or columns, including 13 numerical and three-character variables.

Data Description of Listings

- Name: product name
- Mfr: Manufacturers (A= American Home Food, G=General Mills, K=Kellogs, N=Nabisco, P=Post, Q=Quaker Oats, R=Ralston Purina)
- type: Cereal type (C=cold, H=Hot)
- calories: per serving calories
- protein: grams of protein
- fat: grams of fat
- sodium: milligrams of sodium
- fiber: grams of dietary fiber
- carbo: grams of carbohydrates
- sugar: grams of sugars
- potass: milligrams of potassium
- vitamins: vitamins and minerals (0, 25, or 100)
- shelf: display shelf (1,2,3)
- weight: weight in ounces of one serving
- cups: number of cups in one serving

Module 2-Option 1

- rating: a rating of the cereals

First, load the dataset, then view the statistical summary with the `summary()` function. This function shows each statistical result as a mean value to see an average data point for each variable. Sodium has the highest mean, 159 per serving. Deep end recommends limiting sodium to 1,500 mg daily. Following, calories mean at 106. Mean and median distance is essential for the distribution. The median is much bigger than the mean as the distribution is left-skewed ($\text{median} > \text{mean}$). The median is much smaller than the mean and the right-skewed distribution ($\text{median} < \text{mean}$). The median is equal to or close to 0.7, the normal distribution ($\text{median} = \text{mean}$). The left skew variables are calories (median at 110 > mean at 106) and sodium (median at 180 > mean at 159). The right skew of the variables is potassium (median at 90 < mean at 98), vitamin (median at 25 < mean at 28), and rating (median at 40 < mean at 42). The standard distribution variables are protein, fat, fiber, sugar, shelf, weight, and cups.

A histogram is an excellent representation of category frequency distribution. The histogram plot presents a rectangle for each group of data. Figure 4 has all quantitative columns' results. The `par()` function help to show all variables in one view, makes easy to compare. I used two types of histogram charts. One is the `hist()` function that works with side-by-side plots, and another is the `ggplot2` library. As we talked about distribution on statistical results, also histogram shows the distribution.

A Box plot is an excellent tool to see all the statistical points on graphs like min, median, mode, Q1, Q3, and outliers. Cold cereal Q1 at 100 per serving calorie, and the hot cereal has just one point data at 100. Cold cereal has more data points than hot cereal, and cold cereal has negative and positive side outliers.

Concern

All code is understandable and easy to figure out, but `ggplot` needs clarification. The function gives more errors like the carbon variable error message "Problem while computing aesthetics." I applied for many

Module 2-Option 1

numbers, binwidth, center, and the same error type. I need more detail for the ggplot on why I have this type of error.

Figure 1: Import the Cereal dataset and summary statistics result into R Studio.

The screenshot displays the R Studio interface. The console window shows the following commands and output:

```
> cereal.df <- read.csv("C:/Users/didem/OneDrive/Documents/CSUG Master DA/MIS510-1 Data Mining_4 term/Module 2/MIS510CerealsCSV.csv", header=TRUE)
> #show the first six rows
> head(cereal.df)
```

	name	mfr	type	calories	protein	fat	sodium	fiber	carbo	sugars
1	100%_Bran	N	C	70	4	1	130	10.0	5.0	6
2	100%_Natural_Bran	Q	C	120	3	5	15	2.0	8.0	8
3	All-Bran	K	C	70	4	1	260	9.0	7.0	5
4	All-Bran_with_Extra_Fiber	K	C	50	4	0	140	14.0	8.0	0
5	Almond_Delight	R	C	110	2	2	200	1.0	14.0	8
6	Apple_Cinnamon_Cheerios	G	C	110	2	2	180	1.5	10.5	10

```
> #find summary statistics for each column
> summary(cereal.df)
```

The summary output is as follows:

name		mfr		type		calories				
Length:	77	Length:	77	Length:	77	Min.	:	50.0		
Class	:character	Class	:character	Class	:character	1st Qu.	:	100.0		
Mode	:character	Mode	:character	Mode	:character	Median	:	110.0		
						Mean	:	106.9		
						3rd Qu.	:	110.0		
						Max.	:	160.0		

protein		fat		sodium		fiber		carbo	
Min.	:1.000	Min.	:0.000	Min.	: 0.0	Min.	: 0.000	Min.	: 5.0
1st Qu.	:2.000	1st Qu.	:0.000	1st Qu.	:130.0	1st Qu.	: 1.000	1st Qu.	:12.0
Median	:3.000	Median	:1.000	Median	:180.0	Median	: 2.000	Median	:14.5
Mean	:2.545	Mean	:1.013	Mean	:159.7	Mean	: 2.152	Mean	:14.8
3rd Qu.	:3.000	3rd Qu.	:2.000	3rd Qu.	:210.0	3rd Qu.	: 3.000	3rd Qu.	:17.0
Max.	:6.000	Max.	:5.000	Max.	:320.0	Max.	:14.000	Max.	:23.0

sugars		potass		vitamins		shelf		weight	
Min.	: 0.000	Min.	:15.00	Min.	: 0.00	Min.	:1.000	Min.	:0.50
1st Qu.	: 3.000	1st Qu.	:42.50	1st Qu.	:25.00	1st Qu.	:1.000	1st Qu.	:1.00
Median	: 7.000	Median	:90.00	Median	:25.00	Median	:2.000	Median	:1.00
Mean	: 7.026	Mean	:98.67	Mean	:28.25	Mean	:2.208	Mean	:1.03
3rd Qu.	:11.000	3rd Qu.	:120.00	3rd Qu.	:25.00	3rd Qu.	:3.000	3rd Qu.	:1.00
Max.	:15.000	Max.	:330.00	Max.	:100.00	Max.	:3.000	Max.	:1.50

cups		rating	
Min.	:0.250	Min.	:18.04
1st Qu.	:0.670	1st Qu.	:33.17
Median	:0.750	Median	:40.40
Mean	:0.821	Mean	:42.67
3rd Qu.	:1.000	3rd Qu.	:50.83
Max.	:1.500	Max.	:93.70

The right-hand pane shows the 'Data' tab with a list of loaded datasets: cere..., dmy..., dum..., iris, mmy..., toyo..., and xtot... Each entry shows the number of observations (e.g., 77 obs. for cere...). The 'Values' tab shows the data types: fact... is Factor w/ and samp... is chr [1:5]. The 'Files' and 'Plots' tabs are also visible at the bottom of the right-hand pane.

Figure 2: Display all quantitative variables' statistical results with `data.frame()` function and use the `na.rm=True` function to eliminate missing values in R Studio.

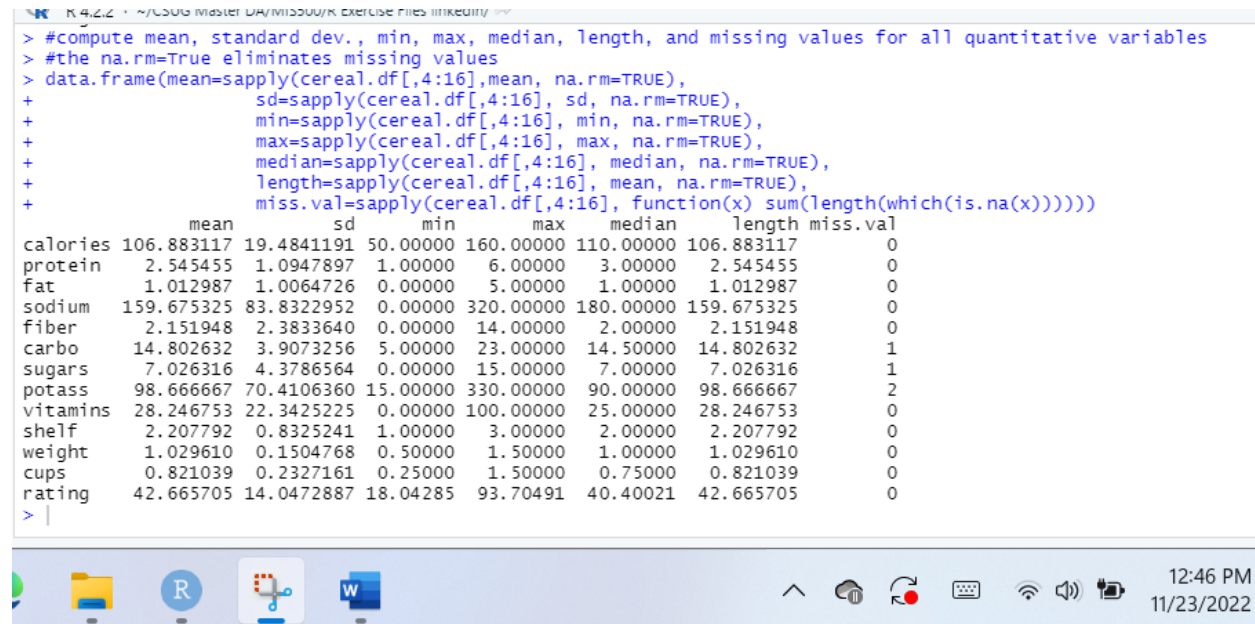


Figure 3: Statistical results of calorie variable in R studio.

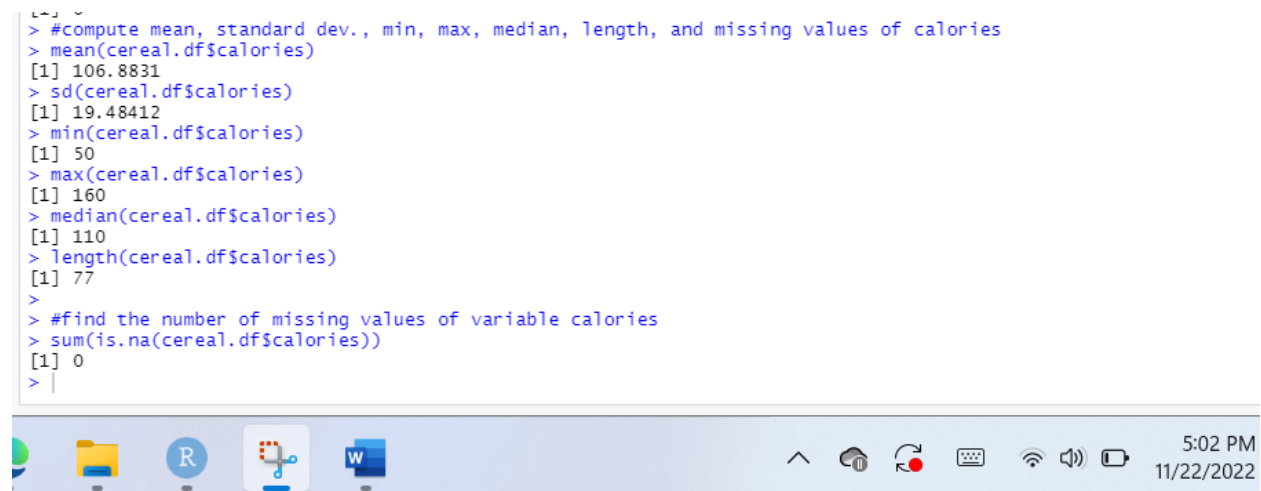


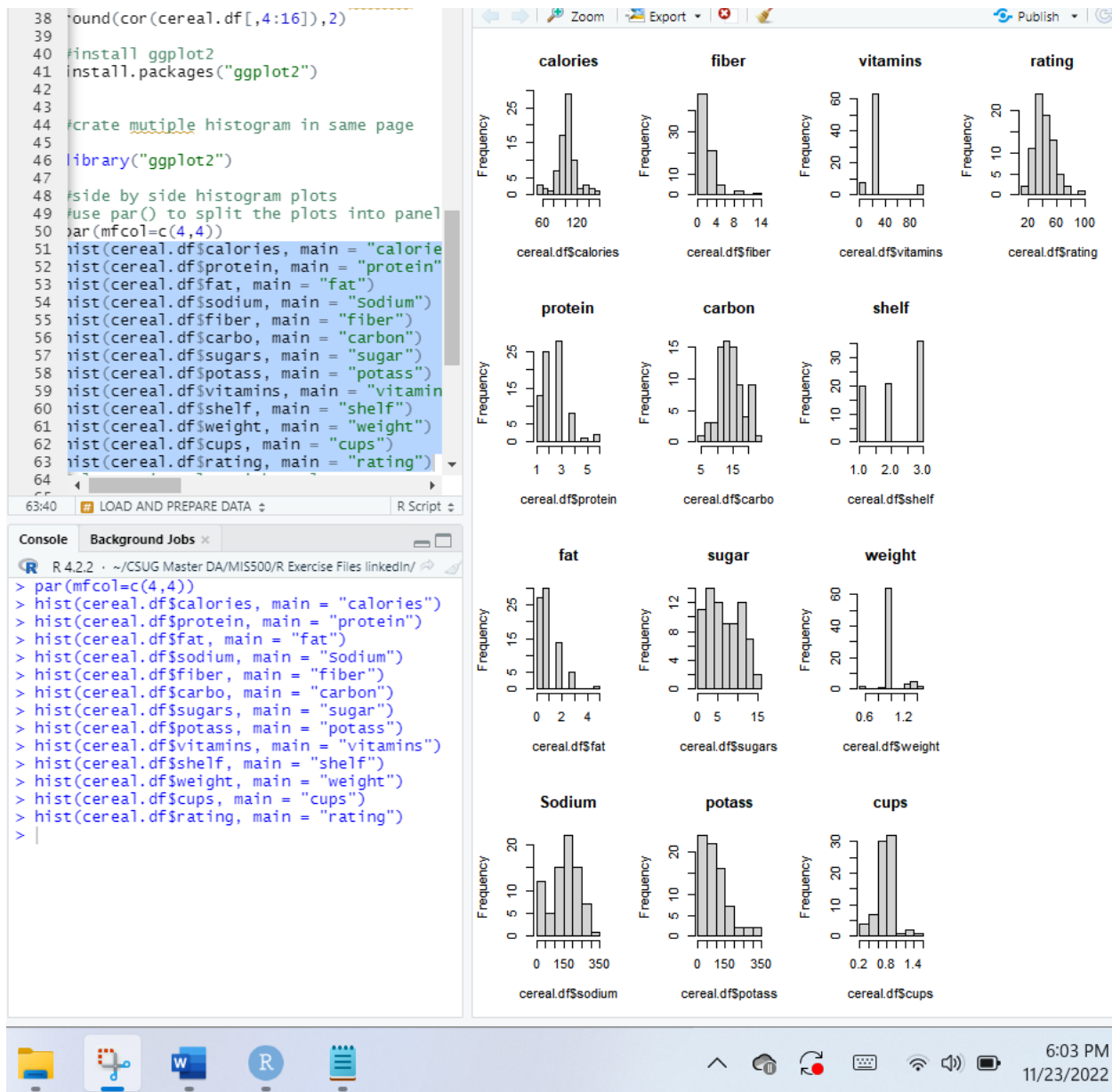
Figure 4: Histogram plot for each quantitative variable side by side with `hist()` function into R Studio.

Figure 5: Alternative way to create histogram plot with `ggplot()` function into R studio.



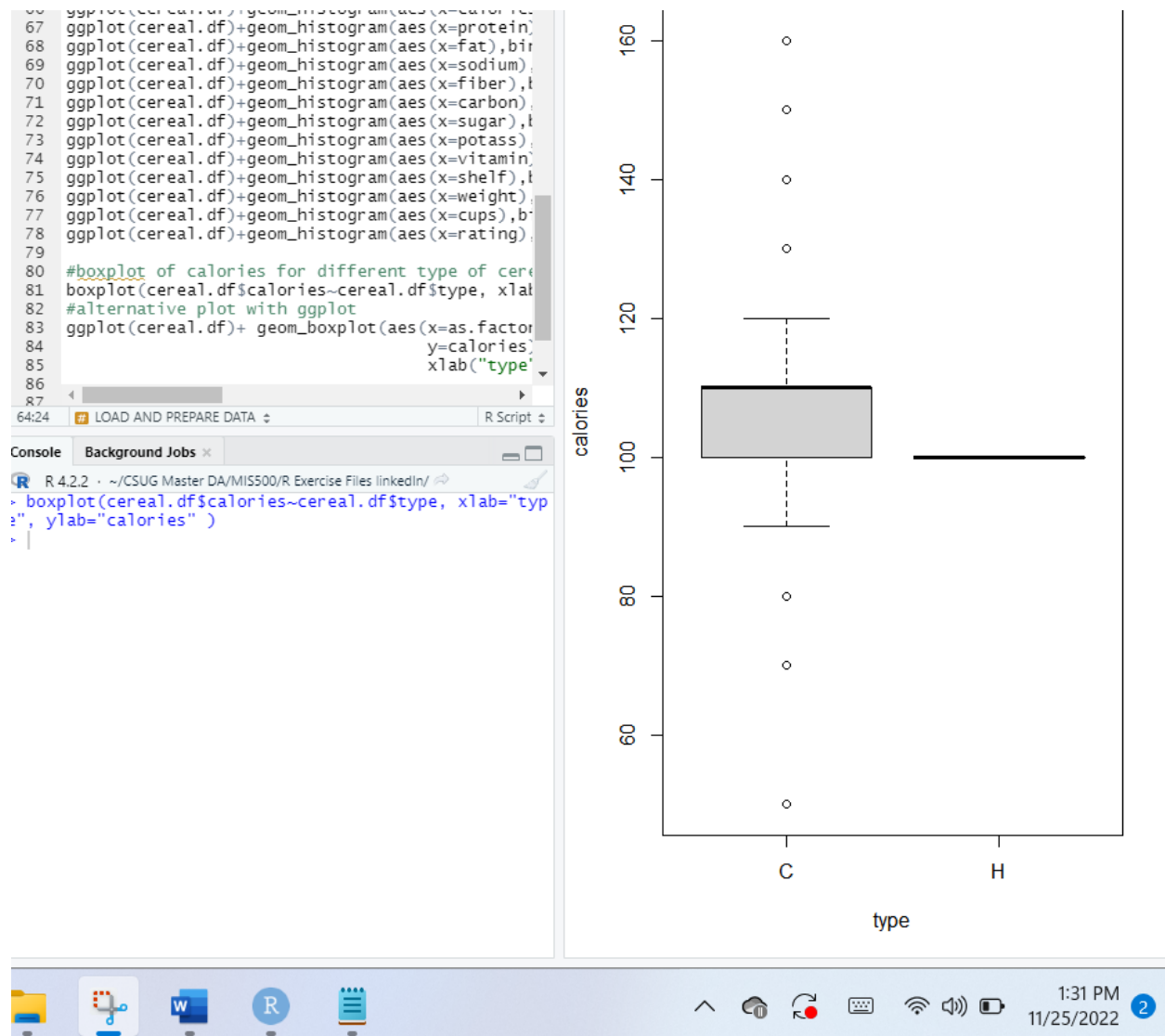
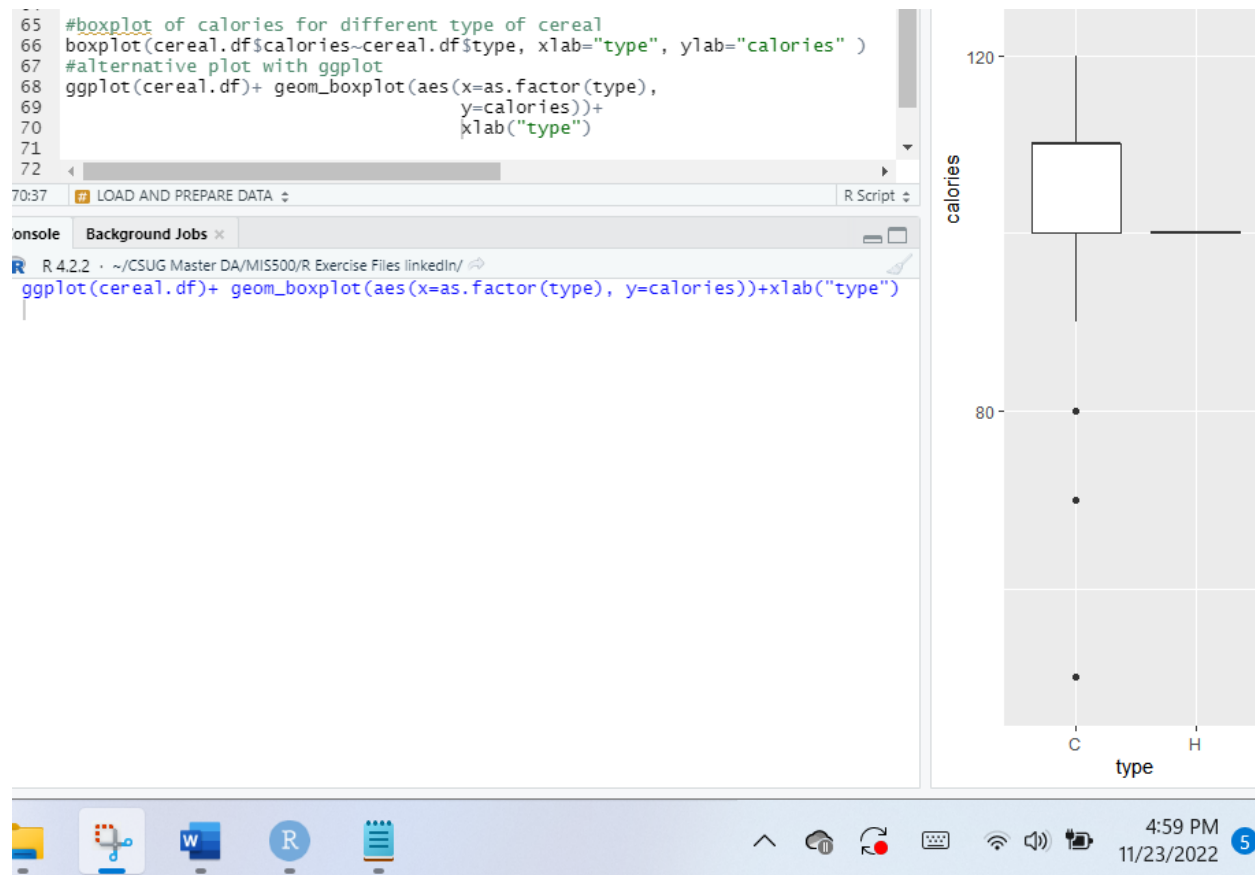
Figure 6: Boxplot for calorie of cereal type with boxplot() function result in R studio.

Figure 7: Alternative way to create boxplot with ggplot() function result in R studio.



Reference

Nutrition data on 80 cereal products, data content by Chris Crawford, 2017

<https://www.heart.org/en/healthy-living/healthy-eating/eat-smart/sodium/how-much-sodium-should-i-eat-per-day>

Written by American Heart Association editorial staff and reviewed by science and medicine advisers. November 1, 2021. <https://www.heart.org/en/healthy-living/healthy-eating/eat-smart/sodium/how-much-sodium-should-i-eat-per-day>

Data Mining for Business Analytics... concepts, Techniques, and Applications in R by Galit Shmueli: Peter C. Bruce: Inbal Yahav: Nitin R. Patel: Kenneth C. Lichtendahl, Jr. Page 56, chapter 3, Figures 3.2 and 3.3. Page 94, chapter 4, Figures 4.3 and 4.4.