

202380-Fall

2023-ITCS-3162-051-Introduction to Data Mining: Project 1

Meelad Doroodchi
09/17/2023

Kaggle Dataset Link:

<https://www.kaggle.com/datasets/rukenmissonnier/age-weight-height-bmi-analysis>

Introduction to Data:

The dataset I am using is named "Age, Weight, Height, BMI Analysis" which has 741 rows of individual records of patients and has 5 columns that include key variables which are: Age, Height(in meters), Weight(in kilograms), Body Mass Index, and Body Mass Index class which includes: "Obese Class 1," "Obese Class 2," "Obese Class 3" "Overweight," "Underweight," "Normal". I am especially interested in this dataset because I believe it is very open-ended but the features allow for the dataset to be visualized and understood in many ways that can be useful in society.

Introduction to Questions:

Some of the questions I would like to solve using the variables would be to understand the relationship between age and body mass index to get an understanding of how age has an impact on someone's body mass index. Along with this, I would like to look at the correlation between weight and body mass index class to be able to get an understanding of the age differences in each BMI class. Looking at different trends with specific ages will be able to help me come to show trends between age groups. The correlation between how much someone weighs and which specific class they are in. Those are just a few of the questions or problems that I would begin to look at with this dataset, I believe that height and body mass index have somewhat of a correlation which allows me to look at a few different relationships to be able to understand the significance of these features on someone's health. A few more questions that I could possibly look at along with these is the distribution of each individual record to show the specific BMI class, BMI index, and the age of each person to be able to look at the overall trend of the dataset.

Additionally, a good thing to look at which BMI in this dataset is the most common and to look at the maximum and minimum to get an understanding of how weight is equivalent to BMI class. This dataset allows me to be able to explore and get an understanding of different trends and connections!

Adv. & Disadv. of Questions:

A few positive benefits that would come from this data is that people would be able to get an understanding of their body types and gain more knowledge knowing who they are. I believe this dataset allows for visualizations and relationships to be shown between age, weight, height, and a person's body and body mass index. This dataset looks at a controlled age range and is a very clean dataset with many open-ended ways to make connections and show that certain relationships can help people gain more understanding!

A few negatives of this dataset include that the age range is from 15-61 which can be both good and bad since it is controlled, but there are many people outside of this age range. As well as this it does not speak about certain health issues that these 741 patients could possibly have that could impact certain things like weight which would ultimately have an impact on body mass index and body mass index class.

Pre-Processing the Data:

Before getting into any visualizations, the data has to be pre-processed. Data cleaning and preprocessing are essential because they ensure data accuracy, improve model performance, handle missing values, outliers, and inconsistencies, standardize data formats, and prepare categorical data for analysis. These tasks create a reliable foundation for data mining, enhancing the quality and effectiveness of the analyses and modeling. In data cleaning a few things that need to be cleaned include the quality of the data. There are many times the data has missing values or the dataset is not fully ready to be used, the programmer has to go back and clean the data before being able to pull any analyses from it. In this project, the dataset that was chosen was quite clean and the implemented data was easy to understand, but to make it simple I changed the weight from kilograms to pounds and the height from meters to feet. I left both of the columns in the dataset so that everyone can be able to view the numbers in both pounds/kilograms and meters/feet.

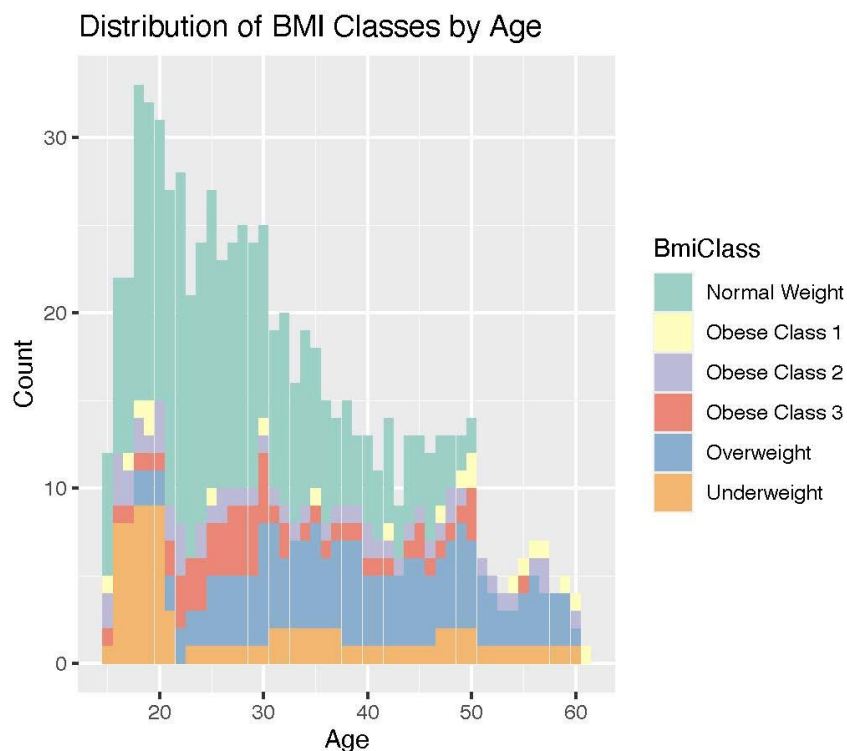
Data Visualization/Story:

The first thing I needed to do before creating any visualizations was to implement any packages I needed into **RStudio**. I downloaded the dplyr and the ggplot2 packages into my console. Using methods in RStudio to create a bar chart, two scatterplots, and a histogram. After looking at some of the questions from above visualizations began to be created to look at relationships. Look at some of the visualizations and connections made below.

Visualization #1: Bar chart of Distribution of BMI classes by Age: Story

Figure 1 below looks at the distribution of BMI classes by age and the number of people in the dataset pertaining to each BMI class. This visualization was a bar chart created in RStudio and this distribution looks at the dataset as a whole and looks at the distribution of age with the BMI class. After looking at the plot some conclusions that were made include the majority of normal weighted people in this dataset ranged from ages 15-50 and not any after the age of 50 which is interesting to see. This looks at the question of age with body mass index which shows that many people over the age of 50 regardless of other health issues will most likely find themselves not at 'normal weight'. This plot shows the BMI class distribution by age and the number of individuals in this dataset in that specific class and age group by color

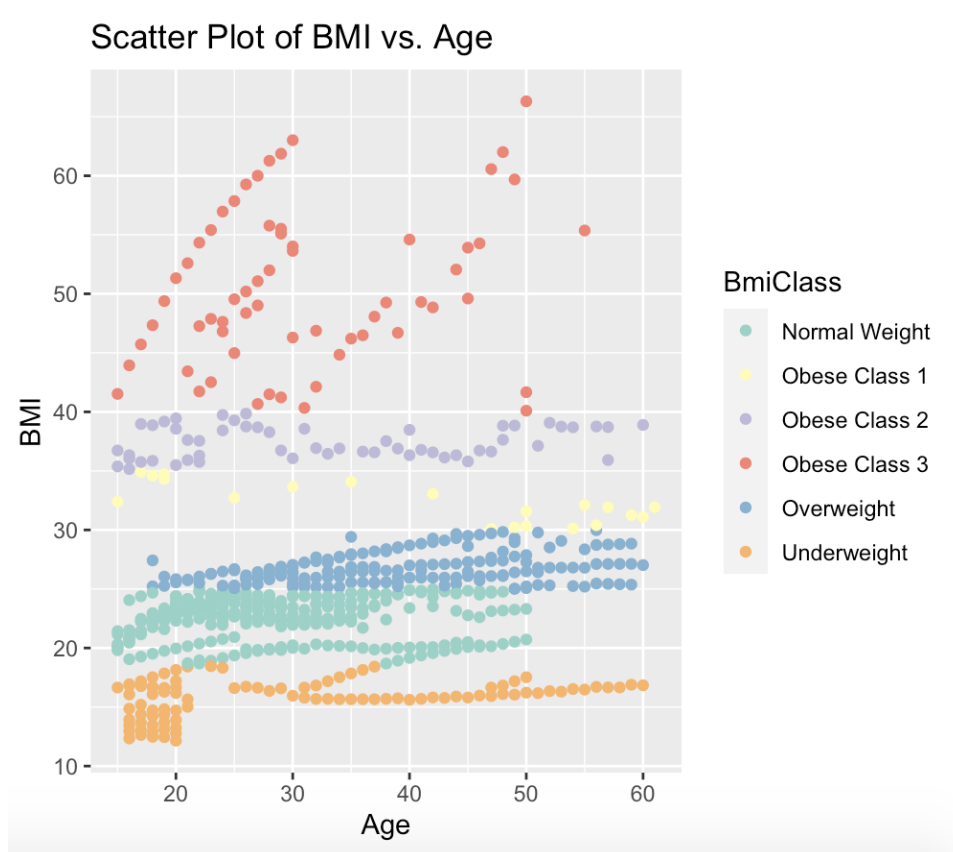
Figure 1. Bar Chart of BMI vs. Age



Visualization #2: Scatterplot of Age vs. BMI and classifying by BMI Class: Story

Figure 2 is a scatter plot of Age vs. BMI while classifying the individual records in this dataset by BMI class as well which answers the following question of what is the specific BMI that is matched with each BMI class and age trends with BMI class and BMI itself. A scatterplot displays the strength, direction, and form of the relationship between two quantitative variables. After looking at this plot some of the answers that were concluded include that the BMI for 'Obese Class 3' is 40+ as it is the heaviest, the BMI for 'Obese Class 2' is 35-39.5, the BMI for 'Obese Class 1' is 30-34.5, the BMI for 'Overweight' is 25-29.5, the BMI for 'Normal Weight' is 19-24.5, and the BMI for 'Underweight' is less than 19. This graph displays different trends to show that this dataset was split in terms of types of patients and records. There are many different chunks of each group significant enough to see on this graph which was interesting to see. It also goes to show that there were not many people past the age of 55 in 'Obese Class 3' which could be due to health.

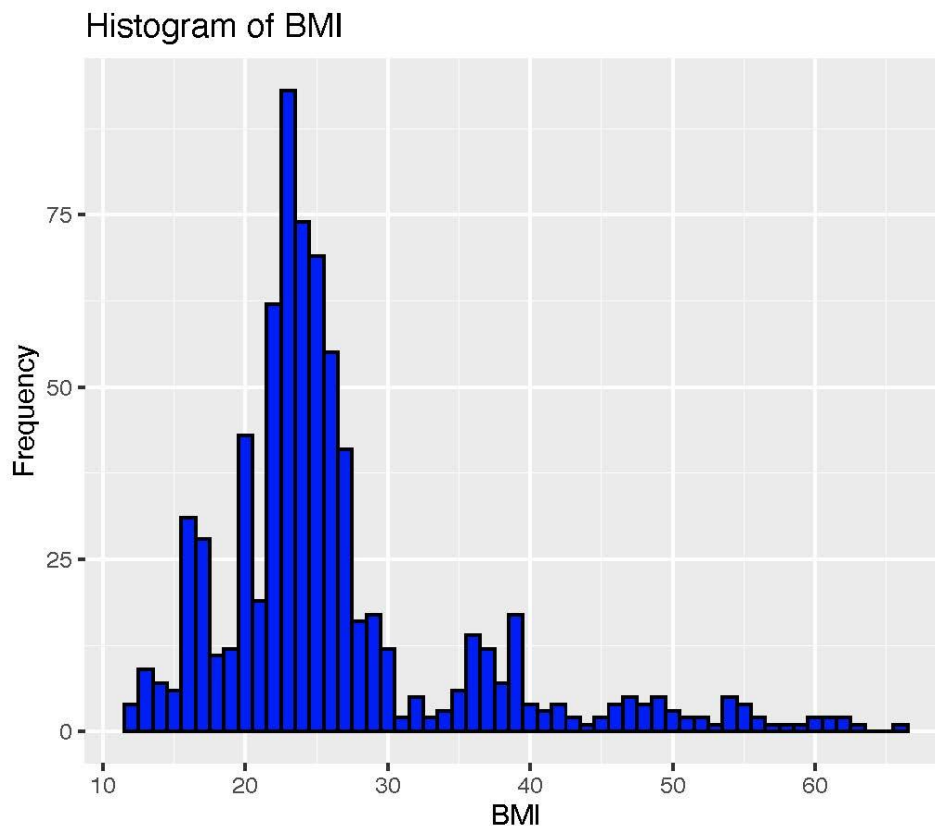
Figure 2. Scatterplot of BMI vs. Age



Visualization #3: Histogram of BMI: Story

Figure 3 below is a Histogram of BMI with frequency or amount of records on this dataset to show the maximum and minimum BMI of this dataset along with the most common BMI. After looking at this dataset some conclusions that were drawn include that the minimum BMI in this dataset is 12 and the maximum is 66 with the most frequent or mode being 23 which is 'Normal Weight' which proves the fact that many individuals in this dataset were of a healthy and normal weight but there were still many different records of people that were severely overweight or underweight according to their BMI's and the conclusions drawn from the previous figures allows to see this dataset as a whole.

Figure 3. BMI Histogram



Visualization #4: Scatterplots of Age vs. Height vs. Weight: Story

Figure 4 below is a scatterplot to show the relationship in this dataset between age, weight, and height. A scatterplot displays the strength, direction, and form of the relationship between two quantitative variables which is why I had to create three different visualizations to show the correlations. Each of these figures is coded by the BMI classification to show the relationship between weight and height, then weight and age, and height and age. From Figure 4 we can see the majority of people under 5'3 are underweight for the most part in this dataset. From Figure 5 we can see that someone who is 100 pounds more than the optimal weight for their sex and height then they are obese class 3. Lastly, in Figure 6 we can see that many people who were under 5'3 were not in these obesity classes for the most part. After looking at all of these we can see the connection that is made with Figure 2 as well with the BMI classes by weight and looking at the connection with height, looking at Obese Class 3 majority of the people are over 40 years old as well.

Figure 4. Scatterplot of Weight vs. Height

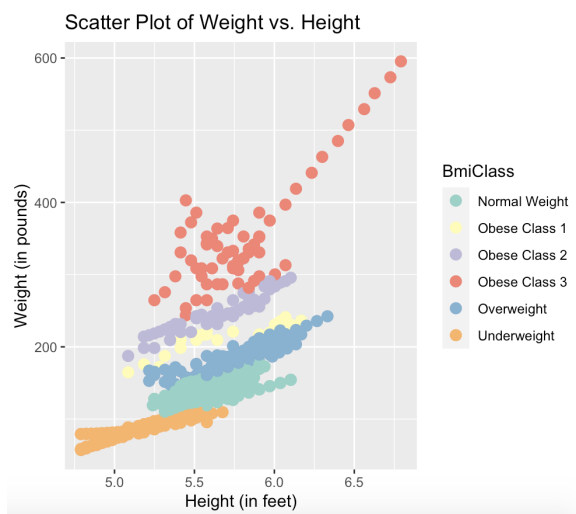


Figure 5. Scatterplot of Weight vs. Age

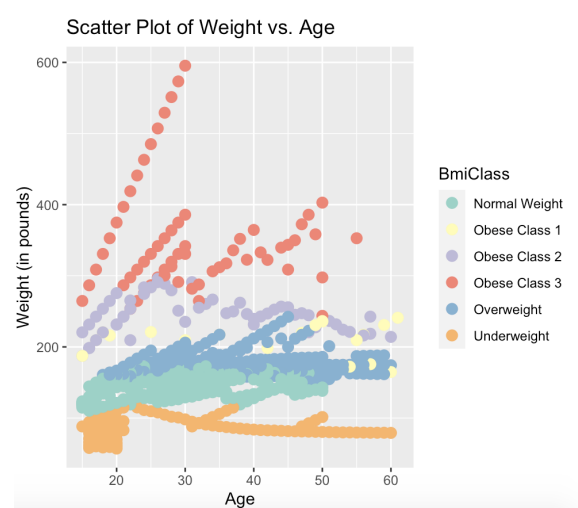
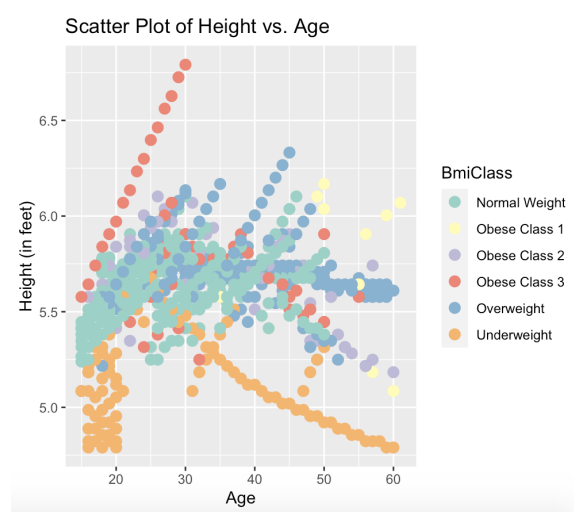


Figure 6. Scatterplot of Height vs. Age



Impact:

The impact of these figures above can be both positive and negative. This dataset with these visualizations is positive because it shows the trends and differences from this dataset which can allow health professionals and even humans themselves to understand their bodies and where they fit into. These four figures above show many things including the bar chart that highlights different age groups and shows how many people past the age of 50 in this dataset were no longer 'Normal Weight' shows that Age can be a deviation in BMI. Including this, these figures specifically figure 2 was able to show the scale of the BMI classes and the threshold for the BMI in each class. With the figure, it is simple to understand where someone can fall. Some negative impacts of this dataset can be generalization, since this dataset is only 741 but over a broad range it can not be generalized for everybody because as stated previously, many different people have underlying health conditions that can affect many different things pertaining to BMI, BMI class, and weight which is not spoken about in this dataset. This study is only looked about from an outer perspective of each person and not certain things each person may have. Overall, these visualizations collectively tell a compelling story about this dataset's characteristics and relationships, and address and show key patterns and trends within the data.

References:

Missonnier, Ruken(1 Sept. 2023). "Age, Weight, Height, BMI Analysis."www.kaggle.com/datasets/rukenmissonnier/age-weight-height-bmi-analysis
. Kaggle. September 17th, 2023

Code: RStudio

Code for Pre-Processing Steps

```
> bmi$Weight_in_Pounds <- bmi$Weight * 2.205  
> bmi$Height_in_Inches <- bmi$Height * 3.28084  
> colnames(bmi)[7] ="Height_in_Feet"
```

Code for Figure 1:

```
> ggplot(data = bmi, aes(x = Age, fill = BmiClass)) +  
+   geom_bar() + # Create a bar chart  
+   labs(title = "Distribution of BMI Classes by Age", x = "Age", y = "Count") +  
+   scale_fill_brewer(palette = "Set3") # Custom color palette (optional)
```

Code for Figure 2

```
> ggplot(data = bmi, aes(x = Age, y = Bmi, color = BmiClass)) +  
+   geom_point() +  
+   labs(title = "Scatter Plot of Age vs. BMI", x = "Age", y = "BMI") +  
+   scale_color_brewer(palette = "Set3") # Choose a color palette
```

Code for Figure 3

```
> ggplot(data = bmi, aes(x = Bmi)) +  
+   geom_histogram(binwidth = 1, fill = "blue", color = "black") +  
+   labs(title = "Histogram of BMI", x = "BMI", y = "Frequency")
```

Code for Figures 4,5,6

Figure 4:

```
> ggplot(data = bmi, aes(x = Height_in_Feet, y = Weight_in_Pounds, color = BmiClass)) +  
+   geom_point(size = 3) +  
+   labs(title = "Scatter Plot of Weight vs. Height",  
+ x = "Height (in feet)", y = "Weight (in pounds)") +  
+   scale_color_brewer(palette = "Set3") # Choose a color palette
```

Figure 5:

```
> ggplot(data = bmi, aes(x = Age, y = Weight_in_Pounds, color = BmiClass)) +  
+   geom_point(size = 3) +  
+   labs(title = "Scatter Plot of Weight vs. Age",  
+ x = "Age", y = "Weight (in pounds)") +  
+   scale_color_brewer(palette = "Set3") # Choose a color palette
```

Figure 6:

```
> ggplot(data = bmi, aes(x = Age, y = Height_in_Feet, color = BmiClass)) +  
+   geom_point(size = 3) +  
+   labs(title = "Scatter Plot of Height vs. Age",  
+ x = "Age", y = "Height (in feet)") +  
+   scale_color_brewer(palette = "Set3") # Choose a color palette
```