# Insert Your Title Here

John Smith
Institution name
Hekla, Iceland
jsmith@affiliation.org

Lars Thørväld
The Thørväld Group
Hekla, Iceland
larst@affiliation.org

Charles Palmer
Palmer Research Laboratories
San Antonio, Texas
cpalmer@prl.com

## ABSTRACT

Data visualization, a rapidly emerging trend, is widely adopted across various sectors. Enterprises and academic institutions leverage the potency of visual representations, such as graphs and charts, to translate intricate data into easily comprehensible formats. Data visualization makes it easy to understand and explain data, helping to make decisions. It is beneficial in any area where you must clearly show significant, complicated information. This research paper provides an extensive overview of data visualization techniques utilizing the R programming language. The paper discusses the data collection process, stressing its importance for reliable visualizations. It then explores the benefits of using R for visualization, emphasizing its versatility and rich package ecosystem. R enables tailored visualizations to meet specific analytical needs. Additionally, the paper covers various visualization techniques in R, such as Pie Diagrams, Bar graphs, Histograms, Box plots, Line plots, Mosaic Plots, Bubble Charts, Violin Plots, and Scatter Diagrams. Each technique is illustrated, offering insights into its strengths and applications. The goal is to equip researchers and practitioners with the knowledge and tools to leverage R for data visualization effectively.

## CCS CONCEPTS

• **Data Visualization**; • **R Language**;

## KEYWORDS

data visualization, Health, R, Correlation

## 1 INTRODUCTION

Today's world is entirely data-driven due to the essential nature of technology. We create new digital footprints every second. The collection and analysis of data have seen an uprise, and it is not going to change. The amount of data we generate every day is a prime reason why the field of Data Science is so important. Data Visualization is a crucial sector within the world of Data Science that explores ways to create visuals from data. Data Visualization allows us to represent data visually compellingly and in an easy-to-understand way. These visuals are not just nice to look at but are essential tools for decision-making and planning in businesses and science [5][6]. We believe performing data visualizations hands-on is the best way to gain knowledge of the various techniques in data visualization. It ensures that we are using these techniques in a practical manner. We also acknowledge that studying existing literature can help us gain further insight into the effectiveness and efficiency of various approaches in Data Visualization. Previous research has mostly explored various types of data visualization techniques and their use cases. We want to build on that and offer more.

In this paper, we study literature that we believe can enrich our understanding of effective communication using data visualization techniques. We then apply most of the applicable data visualization techniques to the attributes of a particular dataset [18] in the R programming language, analyzing every chart and diagram meaningfully throughout the process. Finally, we conclude our research with recommendations for making better Data Visualizations that portray information in engaging ways.

## 2 EARLIER STUDIES

There have been researches that emphasized suggesting guidelines to increase effectiveness in data visualization. Kelleher and Wagener [8] have proposed ten guidelines for creating better data visualizations for scientific journals. Some include: - creating plots that convey the required information, visualizing patterns over information, plotting points that represent density adequately, selecting appropriate color schemes, etc.

Rodríguez et al. [7] survey the historical background of data visualization and its relationship with storytelling. The techniques described by the authors are: - interactive slide show, drill-down story, and martini glass. The authors have also focused on maintaining a balance between communication and engagement. These researches provide valuable insight on conveying information effectively using data visualization tools.

Mohaiminul and Shangzhu [6] discuss the significance of data visualization. They highlight how visualization reveals patterns and trends in data that might be missed in text-based data. The authors also discuss various visualization tools and conclude that effective data visualizations offer key insights into complex datasets.

Emily et al. [2] provide a practical introduction to data visualization using R. The authors highlight the advantages of using R for data visualization, including reproducibility, transparency, and a wide range of fully customizable data visualization options. They explain the 'grammar of graphics' that underlies data visualization using the ggplot package.

Anjali and Rajput [5] underscored the transformation of raw data into meaningful information through efficacious visualization strategies, highlighting the utilization of R for its abilities in statistical analysis and data visualization. They outlined data perception as a tool for improving complex data for more precise communication and faster decision-making. This paper represents the significance and use of data visualization in extricating significant understanding using R.

Dengyun et al. [3] explored innovative approaches to data visualization in civil aviation. The researchers emphasized the crucial role of advanced visualization techniques in improving safety measures and operational efficiency. They suggested coordinating these approaches into standard practices to deal with the undeniably complicated datasets in aeronautics (ICASIT, 2023). The overview focused on the latest trends and set a forward-looking plan for taking on these technologies on a broader scale.

Yavuz S. Taspinar et al. used a machine learning approach to determine sleep health status, distinguishing between sleep apnea, insomnia, and normal sleep [1]. Using a dataset with 374 rows and 13 features, they applied Random Forest, SVM, Logistic Regression, and k-NN. The Random Forest model achieved the highest classification success rate of 91.66%. Body mass index was found to have the most significant impact on diagnosing sleep disorders, as revealed by box plot and heatmap correlation analysis.

Stephen R. Halfway et al. [4] proposed ten principles for effective data visualization, emphasizing the importance of projecting key statistics in histograms and using different colors to distinguish between data and display. The paper also stresses the need for precise definitions of infographic data, captions, and visualization nuances.

## 3 METHODOLOGY

In this section, the data collection process is discussed first. Later, the benefits of using R as data visualization are explained, and finally, several data visualization techniques are illustrated with the outcome of each diagram.

### 3.1 Data Collection

The dataset for this research was collected from Kaggle, a popular online platform for data science and machine learning. The "Sleep Health and Lifestyle Dataset" dataset provides a comprehensive view of various health and lifestyle factors and their impact on sleep [18].

The dataset comprises several attributes, each representing a specific aspect of an individual's health, lifestyle, or demographic information. The following are the attributes included in the dataset:

- `Person ID`: An identifier for each individual [18].
- `Gender`: The person's gender (Male/Female) [18].
- `Age`: The age of the person in years [18].
- `Occupation`: The occupation or profession of the person [18].
- `Sleep Duration (hours)`: The number of hours the person sleeps daily [18].
- `Quality of Sleep (scale: 1-10)`: A subjective sleep quality rating ranging from 1 to 10 [18].

- `Physical Activity Level (minutes/day)`: The number of minutes the person engages in physical activity daily [18].
- `Stress Level (scale: 1-10)`: A subjective rating of the stress level experienced by the person, ranging from 1 to 10 [18].
- `BMI Category`: The BMI category of the person (e.g., Normal, Obese, Overweight) [18].
- `Blood Pressure`: The blood pressure measurement of the person [18].
- `Heart Rate (bpm)`: The person's resting heart rate in beats per minute [18].
- `Daily Steps`: The number of steps the person takes daily [18].
- `Sleep Disorder`: The presence or absence of sleep disorder in the person (None, Insomnia, Sleep Apnea) [18].

This dataset allows for a wide range of analyses, from exploring the relationship between sleep disorder and occupation to investigating the impact of stress levels on sleep quality.

### 3.2 Data Visualization Using R

Data Visualization is the process of creating graphical representations of information. This process helps the presenter communicate data in a way that's easy for the viewer to interpret and draw conclusions [10]. It involves the use of visual elements like charts, graphs, and maps, which provide an accessible way to see and understand trends, outliers, and patterns in data [10][11][12][13].

R, a popular programming language for statistical computing, offers robust capabilities for data visualization [14][15][16]. It provides a broad collection of visualization libraries and extensive online guidance on their usage [16]. We can easily customize our data visualization through R by changing axes, fonts, legends, annotations, and labels [16]. One advantage of using R is its potential benefits to reproducibility and transparency.

Several common data visualization techniques can be implemented using R. These include:

- `Pie Diagram`: Ideal for illustrating proportions or part-to-whole comparisons [17].
- `Bar Graph`: Perfect for comparing quantities among different groups [17].
- `Histogram`: Used to show the distribution of a dataset [17].
- `Box plot`: Ideal for showing a dataset's range and other characteristics [17].

These techniques can be implemented using various packages in R, such as ggplot2, which is part of the tidyverse collection of packages [14]. The ggplot2 package allows users to create various visualizations, from basic plots like bar charts and histograms to more complex visualizations.

### 3.3 Pie Diagram

Pie diagrams, also known as pie charts, are widely used graphical representations that display data in a circular format [5]. Each 'slice' of the pie represents a proportion or percentage of the whole dataset, with the size of each slice corresponding to its relative magnitude.

The pie diagram of the BMI [figure 1] represents the distribution of different Body Mass Index (BMI) categories in a population.
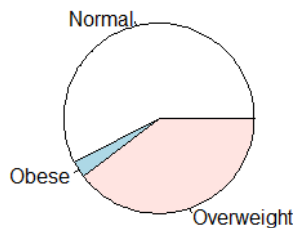
**Pie Diagram of BMI**



**Figure 1: Pie Diagram to Demonstrate the Distribution of BMI Category.**

Most of the population falls within the 'Normal' BMI range (216 instances), indicating a healthy weight status. However, a significant portion is categorized as 'Overweight' (148 instances), suggesting a trend towards weight issues in this population. The 'Obese' category is the most minor (10 instances), but it is still a concern. The pie chart effectively visualizes these distributions, highlighting the prevalence of the 'Normal' and 'Overweight' categories and the relatively low frequency of the 'Obese' category.

## 3.4 Bar Graph

A bar graph, also known as a bar chart, is a type of graph that represents a categorical variable with columns plotted vertically or horizontally. A categorical variable has two or more categories with no inherent order [22]. Bar graphs are useful for comparing data across different categories, with the length of the bars representing the number of observations for each category.
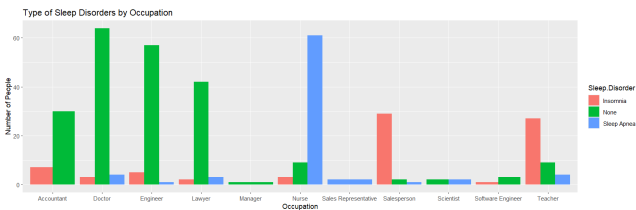


**Figure 2: Bar Graph Demonstrating Type of Sleep Disorder by Occupation.**

The bar graph shown in [figure 2] represents the number of individuals with sleep disorders categorized into three groups: insomnia, sleep apnea, and none, based on occupation. The graph indicates that teachers and salespersons are more likely to suffer from insomnia. Additionally, it shows that nurses have the highest tendency to have sleep apnea, while accountants, doctors, engineers, and lawyers are the least likely to have any sleep disorder.

## 3.5 Histogram

A histogram is a graphical representation of the distribution of numerical data [5]. It consists of a series of contiguous bars, each representing the frequency or relative frequency of values falling within a specified range or 'bin' of the data. Histograms are widely used in data analysis to visualize a dataset's shape, center, and spread.
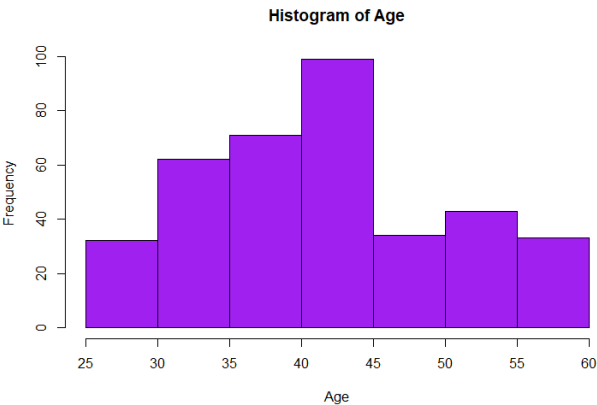


**Figure 3: Histogram to Demonstrate the Frequency of Age.**

The histogram of the Age [figure 3] has a unimodal shape with a peak around the 40-45 age range, indicating this is the most common age group. The center or mode of the distribution is around the mid-40s, and the ages span from 25 to 60, with most data points falling between approximately 35 and 55 years old. The pattern of this histogram suggests that individuals in their early to mid-40s are the most frequent, while those in their late 50s are less frequent. There do not appear to be any significant outliers. The histogram is slightly right-skewed, indicating a longer tail towards older ages. However, it is not strongly skewed as many individuals are still in older age categories.

## 3.6 Box Plot

A box plot, also known as a box-and-whisker plot, is a visual representation that provides a clear summary of the distribution of numerical data. The box in the plot represents 50% of the data, with the lower and upper edges indicating the first and third quartiles, respectively. A line inside the box marks the median value. The 'whiskers,' or lines that extend from each end of the box, show the range of remaining data, excluding outliers. Outliers are values that lie outside this range and are plotted individually as dots.[19].

The box plot [figure 4] shows heart rate data with a median of 70, mean of 70.17, and mode of 68. The distribution is relatively symmetrical, with some high-end outliers. These outliers suggest occasional unusually high heart rates. Despite these outliers, the mean, median, and mode proximity indicate a lack of significant skewness in the data. The data is reasonably consistent, with a few exceptions, such as higher rates. The most common heart rate (mode) falls within the middle 50% of data, further supporting the lack of skewness. The box plot summarizes the heart rate data, highlighting its central tendencies, spread, and potential outliers.
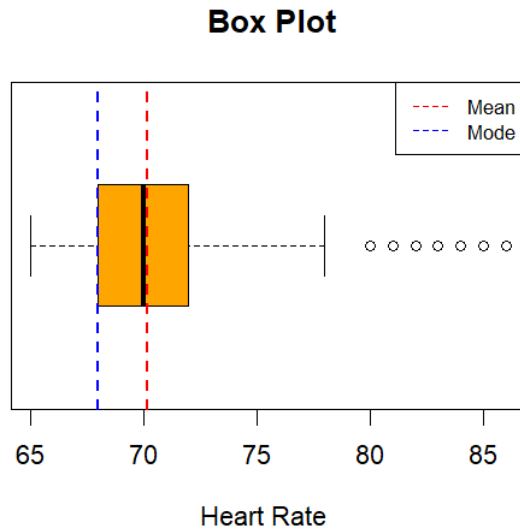
## Box Plot



**Figure 4: Box Plot to Demonstrate Heart Rate Data.**

### 3.7 Line Plot

Line plots are essential tools in data visualization, mainly for displaying trends and patterns over continuous variables such as time or age. They consist of points connected by straight lines, each representing data value at a specific location on the x-axis [5]. Line plots are commonly used to visualize relationships, changes, and trends in data over time or across different conditions.
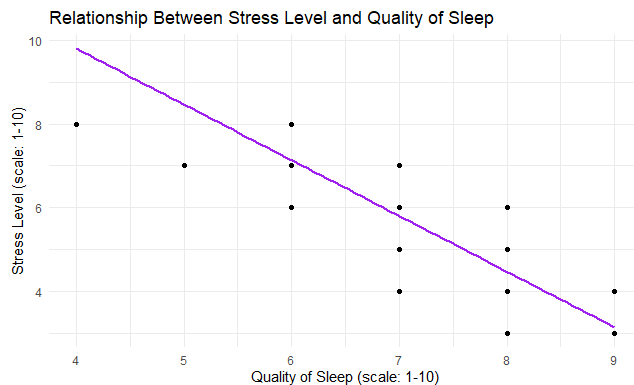


**Figure 5: Line Plot to Demonstrate the Relationship between Stress Level and Quality of Sleep.**

The line plot [figure 5] depicts the relationship between stress level and quality of sleep. It shows a clear negative correlation, indicating that as the quality of sleep improves (increases on the x-axis), the stress level decreases (decreases on the y-axis). This trend suggests that better sleep is associated with lower stress levels. The data points are closely aligned with the trendline, indicating a strong relationship between these variables. The downward-sloping

pattern reinforces the inverse relationship between sleep quality and stress level. This insight could be valuable in understanding and improving individual health and well-being.

### 3.8 Mosaic Plot

A mosaic plot is a stacked bar chart to visualize the relationship between two categorical variables. In the plot, the width of the columns represents the number of observations in each category of the first variable along the horizontal axis. Conversely, the height of the bars corresponds to the number of observations in the second variable within each category of the first variable. Mosaic plots help compare different groups and illustrate relationships between variables clearly and concisely.[20].
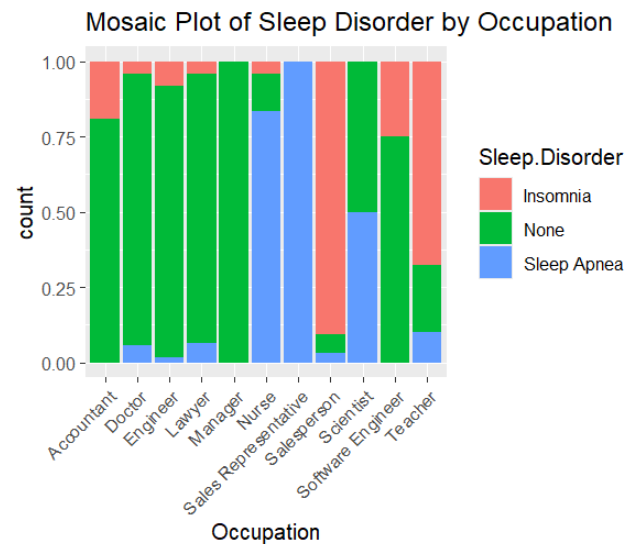


**Figure 6: Mosaic Plot to Demonstrate the Relationship between Occupation and Sleep Disorder.**

The Mosaic plot [figure 6] depicts the relationship between occupation and sleep disorder attributes using different colors such as green, sky blue, and light pink to indicate the rates of "none," "sleep apnea," and "insomnia." In the plot, the x-axis lists various occupations, and the y-axis depicts the count rate for each sleep disorder. For instance, the plot shows that approximately 5% of nurses do not develop any sleep disorder, while nearly 80% experience sleep apnea and roughly 15% experience insomnia. This relationship suggests a possible correlation between sleep disorder and occupation, indicating that the type of sleep disorder may depend on one's occupation.

### 3.9 Bubble Chart

Bubble charts are a type of data visualization that simultaneously represents three data dimensions: two numerical variables represented on the x and y axes and a third numerical variable represented by the size of the bubbles. In a bubble chart, each data point is represented by a circle (or bubble), where the position of the bubble on the x and y axes corresponds to the values of the two

variables, and the size of the bubble represents the magnitude of the third variable [21].
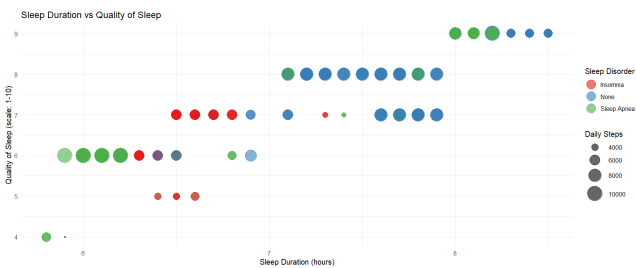


**Figure 7: Bubble Chart to Demonstrate the Connection between Sleep Duration, Sleep Quality, and Daily Steps.**

The connection between sleep duration, sleep quality, and day-to-day step count for people with different sleep problems is outwardly addressed by the Bubble Chart Diagram [figure 7]. The graph indicates that individuals without sleep problems (blue bubbles) regularly have good sleep duration. The length of sleep for individuals with sleep apnea (green bubbles) changes, yet strangely, their sleep quality is not generally inferior; all things considered, it is typically disseminated through the center to the higher finish of the range. Moreover, the predictable proof is absent across the information, focusing on the connection between the number of daily steps (bubble size) and the quality or length of sleep. The information shows fluctuation, particularly among the people who have insomnia (red bubbles) and display a variety in the duration and nature of their sleep. Generally, the effect of everyday steps on sleep quality and duration is more subtle, even though the outline shows a few expansive patterns and connections between sleep disorders, sleep measurements, and active work.

### 3.10 Violin Plot

A violin Plot is a technique to show the conveyance of numerical data of various variables. A violin plot is a hybrid of a box plot and a kernel density plot, which shows peaks in the data. Dissimilar to a box plot that can show statistical insights, violin plots portray summary statistics and the frequency of every variable [9].
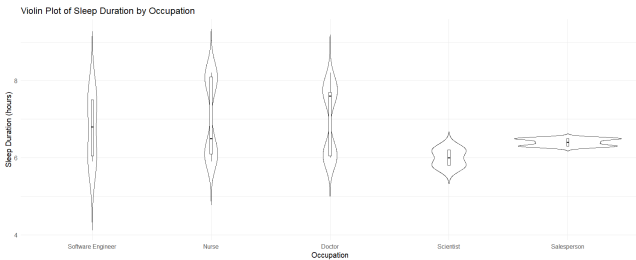


**Figure 8: Violin Plot to Demonstrate the Connection between Occupation and Sleep Duration.**

The Violin plot [figure 8] visualizes the connection between occupation and sleep duration. Scientists experience the lowest

median sleep duration at 6 hours, showing expected difficulties in keeping up with adequate rest in this field. Salespersons follow intimately, with the second-least median sleep duration of around 6.5 hours. While most occupations display some level of balance, Scientists and Salespeople show a more articulated right skewness. This demonstrates that more people in these two professions are confronted with more limited rest lengths than in other professions.

### 3.11 Scatter Diagram

The scatter diagram graphs pairs of numerical data, with one variable on each axis, to look for a relationship between them. The points will fall along a line or curve if the variables are correlated. The better the correlation, the tighter the points hug the line [23]. Scatter diagrams are handy for quickly visualizing correlations among attributes.
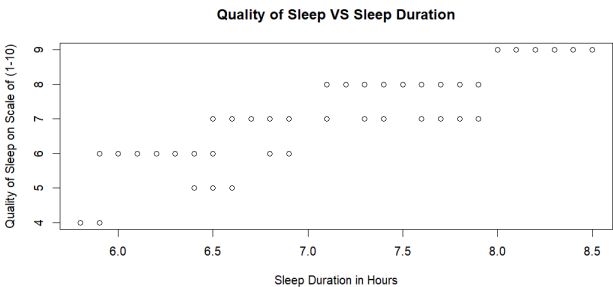


**Figure 9: Scatter Diagram to Demonstrate the Correlation between Quality of Sleep and Sleep Duration.**

The scatter diagram [figure 9] shows Sleep Duration (in hours) plotted against Quality of Sleep (on a scale of 1 − 10). According to the diagram, it can be observed that the attributes of Sleep Duration and Quality of Sleep have a positive correlation. This is denoted by the linear increase in the Quality of Sleep as the Sleep Duration increases. It is also observed that the highest quality of sleep can be obtained if the sleep duration is between 8 and 8.5 hours, while the lowest quality of sleep occurs when the sleep duration is less than 6 hours.

## 4 CHALLENGES AND LIMITATIONS

Our initial challenge was selecting a suitable and reliable dataset for our analysis. We carefully selected potential datasets, prioritizing data quality and relevance to our research question. Once we had ensured a strong dataset, we focused on determining the most appropriate univariate and multivariate analysis techniques. This involved considering the nature of our data and the potential insights each approach could offer. Throughout the process, we carefully searched for reliable and authentic sources of information, prioritizing academic journals, reputable institutions, and peer-reviewed studies to ensure the accuracy and reliability of our research.

We mainly focus on descriptive visualization and do not delve into inferential analysis or explore multiple datasets. We may use a pie chart to show how we spend our time, but it would not show how

our mood changes throughout the day. A bar chart could compare our favorite books, but it would not capture the complex emotions they evoke. Each visualization method has its strengths but also blind spots—Histograms group things together, losing individual stories. Box plots can be skewed by extreme cases. Scatter plots struggle with crowds, and bubble charts get messy with too many dimensions. Even the elegant violin plot can leave some scratching their heads.

## 5 DISCUSSION

This research study delves into various approaches to data visualization utilizing the R programming language, starting with a comprehensive overview of data collection methods and highlighting the advantages of employing R for visualization tasks. Through the practical application of visualization techniques such as Pie Diagrams, Bar Graphs, Histograms, Box Plots, Line Plots, Mosaic Plots, Bubble Charts, Violin Plots, and Scatter Diagrams on a specific dataset, the study offers detailed insights into each method. It enriches the existing body of knowledge by reviewing previous research in the field and demonstrating the effective use of R to communicate complex information visually. Integrating theoretical concepts with practical applications gives researchers and practitioners the tools and insights to leverage R for data visualization effectively.

In the practical realm, these visualization techniques are applied across various sectors to facilitate better decision-making and strategic planning. For instance, pie diagrams and bar graphs are extensively used in business and education to represent financial distributions and academic performance comparisons. Histograms and Box Plots are pivotal in market research and quality control for analyzing customer demographics and product variability. Line Plots track trends over time in business sales and educational enrollments, while Mosaic and Bubble Charts offer advanced segmentation and multi-dimensional analysis in business contexts. Violin Plots and Scatter Diagrams are utilized in educational assessments and human resource management to explore data distributions and correlations between variables like employee satisfaction and productivity. By choosing appropriate visualization techniques for specific data types, this study enhances the descriptive capabilities of data representation. It contributes to predictive and prescriptive analytics, significantly influencing policy-making, business strategy development, and educational advancements.

## 6 CONCLUSION

In this paper, we have explored many useful data visualization techniques using the R programming language. Data Visualization is an essential tool for conveying information meaningfully. Also, we must choose the correct plots and diagrams for different attributes in our datasets. An adequate visualization technique can ensure we interpret the right information from our data and determine the correct relationships among data attributes. We then applied and analyzed many visualization techniques for a chosen dataset and described the types of visualization techniques according to the type of attributes. A thorough literature review of the existing research papers focusing on data visualization techniques was also accomplished. The reviews have brought deep insight into the

effective and efficient ways to visualize data to make them more accessible and interactive. In conclusion, this paper is a strong foundation for future researchers who want to delve into the world of data science and data visualization.