# ASSIGNMENT 1

- **Shaurya Singh Srinet**

**DATASET DETAILS-** Sample Dataset taken from google

sample-data-10mins.
xlsx

## 1. LOADING AND OVERVIEW OF THE DATA:

## 2. DEFINING NEW DATA POINTS:



## 3. COMPARISONS:

**Screenshot 1 — RStudio: assignment-1 - master**

Data viewer tab: sample_data_10mins — Showing 44 to 53 of 1,094 entries, 6 total columns

| | Sales Person | Country | Product | Date | Amount | Boxes Sh |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | |
| 52 | Brien Boise | Australia | Fruit & Nut Bars | 2022-06-06 | 8575 | |

Environment:

Data
- country_sales — 6 obs. of 4 variables
- data.df — 1094 obs. of 7 variables
- sample_data_1... — 1094 obs. of 6 variables

Files: g Data Visualization > Assignments > Assignment 1 > Dataset > assignment-1

| | Name | Size | Modified |
|---|---|---|---|
| | .. | | |
| | assignment-1.Rproj | 218 B | Jun 16, 2024, 9:55 PM |

Console:

```
R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/
> summary(data.df$Amount_per_Box)
    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
   0.014  14.593  36.142 105.720  79.160 4291.000
> # Comparisons: Country-wise sales
> country_sales <- data.df %>%
+     group_by(Country) %>%
+     summarise(Total_Amount = sum(Amount),
+               Average_Amount_per_Box = mean(Amount_per_Box),
+               Total_Boxes = sum(`Boxes Shipped`))
>
> print(country_sales)
# A tibble: 6 × 4
  Country    Total_Amount Average_Amount_per_Box Total_Boxes
  <chr>             <dbl>                  <dbl>       <dbl>
1 Australia       1137367                   95.4       32647
2 Canada           962899                  121.        31221
3 India           1045800                   92.9       29470
4 New Zealand      950418                  118.        26580
5 UK              1051792                   90.8       30265
6 USA             1035349                  119.        26824
>
```

10:15 PM 16-06-2024

---



**Screenshot 2 — RStudio: assignment-1 - master**

Data viewer tab: sample_data_10mins — Showing 44 to 52 of 1,094 entries, 6 total columns

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Environment:

Data
- country_sales — 6 obs. of 4 variables
- data.df — 1094 obs. of 7 variables
- product_sales — 22 obs. of 4 variables
- salesperson_sales — 25 obs. of 4 variables
- sample_data_10mins — 1094 obs. of 6 variables

Files: Shaurya > Desktop > SEM 7 SUBJECTS > Big Data Visualization > Assignments > Assignment 1 > Dataset > assignment-1

| | Name | Size | Modified |
|---|---|---|---|
| | .. | | |
| | assignment-1.Rproj | 218 B | Jun 16, 2024, 9:55 PM |

Console:

```
R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/
> # Contrasts: Salesperson-wise sales
> salesperson_sales <- data.df %>%
+     group_by(`Sales Person`) %>%
+     summarise(Total_Amount = sum(Amount),
+               Average_Amount_per_Box = mean(Amount_per_Box),
+               Total_Boxes = sum(`Boxes Shipped`))
>
> print(salesperson_sales)
# A tibble: 25 × 4
   `Sales Person`     Total_Amount Average_Amount_per_Box Total_Boxes
   <chr>                     <dbl>                  <dbl>       <dbl>
 1 Andria Kimpton           201747                  106.         6448
 2 Barr Faughny             258713                  178.         6366
 3 Beverie Moffet           278922                   80.5        9214
 4 Brien Boise              312816                   75.8        8102
 5 Camilla Castle           196616                   89.6        5374
 6 Ches Bonnell             320901                  106.         7522
 7 Curtice Advani           216461                  114.         7074
 8 Dennison Crosswaite      291669                   76.8        8767
 9 Dotty Strutley           190624                  112.         6853
10 Gigi Bohling             232666                   93.3        6303
# i 15 more rows
# i Use `print(n = ...)` to see more rows
>
```

10:16 PM 16-06-2024

# 4. CONTRASTS:



# 5. IDENTIFYING TENDENCIES:

# 6. DISPERSIONS:



# 7. PLOTTING:

sample_data_10mins

Filter

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madeiene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

Environment  History  Connections  Git  Tutorial

Import Dataset · 251 MiB · 

R · Global Environment

Data

| | | |
|---|---|---|
| country_sales | 6 obs. of 4 variables | |
| data.df | 1094 obs. of 7 variables | |
| dispersions | 1 obs. of 6 variables | |
| product_sales | 22 obs. of 4 variables | |
| salesperson_sales | 25 obs. of 4 variables | |
| sample_data_10mins | 1094 obs. of 6 variables | |
| tendencies | 1 obs. of 5 variables | |

Console  Terminal  Background Jobs

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```
+          IQR_Amount = IQR(Amount),
+          Variance_Amount_per_Box = var(Amount_per_Box),
+          SD_Amount_per_Box = sd(Amount_per_Box),
+          IQR_Amount_per_Box = IQR(Amount_per_Box))
> print(dispersions)
# A tibble: 1 × 6
  Variance_Amount SD_Amount IQR_Amount Variance_Amount_per_Box SD_Amount_per_Box
            <dbl>     <dbl>      <dbl>                   <dbl>             <dbl>
1       16830030.     4102.      5637.                  77651.              279.
# i 1 more variable: IQR_Amount_per_Box <dbl>
> # Plotting: Country-wise total amount
> ggplot(country_sales, aes(x = Country, y = Total_Amount)) +
+     geom_bar(stat = "identity") +
+     theme_minimal() +
+     labs(title = "Country-wise Total Amount",
+          x = "Country", y = "Total Amount")
> # Plotting: Product-wise total amount
> ggplot(product_sales, aes(x = Product, y = Total_Amount)) +
+     geom_bar(stat = "identity") +
+     theme_minimal() +
+     labs(title = "Product-wise Total Amount",
+          x = "Product", y = "Total Amount")
>
```

Files  Plots  Packages  Help  Viewer  Presentation

Zoom  Export  Publish



Product-wise Total Amount (bar chart with Total Amount on y-axis, Product on x-axis)

---

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madeiene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

Environment  History  Connections  Git  Tutorial

Import Dataset · 252 MiB · 

R · Global Environment

Data

| | | |
|---|---|---|
| country_sales | 6 obs. of 4 variables | |
| data.df | 1094 obs. of 7 variables | |
| dispersions | 1 obs. of 6 variables | |
| product_sales | 22 obs. of 4 variables | |
| salesperson_sales | 25 obs. of 4 variables | |
| sample_data_10mins | 1094 obs. of 6 variables | |
| tendencies | 1 obs. of 5 variables | |

Console  Terminal  Background Jobs

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```
# A tibble: 1 × 6
  Variance_Amount SD_Amount IQR_Amount Variance_Amount_per_Box SD_Amount_per_Box
            <dbl>     <dbl>      <dbl>                   <dbl>             <dbl>
1       16830030.     4102.      5637.                  77651.              279.
# i 1 more variable: IQR_Amount_per_Box <dbl>
> # Plotting: Country-wise total amount
> ggplot(country_sales, aes(x = Country, y = Total_Amount)) +
+     geom_bar(stat = "identity") +
+     theme_minimal() +
+     labs(title = "Country-wise Total Amount",
+          x = "Country", y = "Total Amount")
> # Plotting: Product-wise total amount
> ggplot(product_sales, aes(x = Product, y = Total_Amount)) +
+     geom_bar(stat = "identity") +
+     theme_minimal() +
+     labs(title = "Product-wise Total Amount",
+          x = "Product", y = "Total Amount")
> # Plotting: Salesperson-wise total amount
> ggplot(salesperson_sales, aes(x = `Sales Person`, y = Total_Amount)) +
+     geom_bar(stat = "identity") +
+     theme_minimal() +
+     labs(title = "Salesperson-wise Total Amount",
+          x = "Sales Person", y = "Total Amount")
>
```

Files  Plots  Packages  Help  Viewer  Presentation

Zoom  Export  Publish



Salesperson-wise Total Amount (bar chart with Total Amount on y-axis, Sales Person on x-axis)

# Screenshot 1

**assignment-1 - master - RStudio**

File Edit Code View Plots Session Build Debug Profile Tools Help

sample_data_10mins

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

**Environment / History / Connections / Git / Tutorial**

Data
- country_sales — 6 obs. of 4 variables
- data.df — 1094 obs. of 7 variables
- dispersions — 1 obs. of 6 variables
- product_sales — 22 obs. of 4 variables
- salesperson_sales — 25 obs. of 4 variables
- sample_data_10mins — 1094 obs. of 6 variables
- tendencies — 1 obs. of 5 variables

**Console / Terminal / Background Jobs**

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```
+   geom_bar(stat = "identity") +
+   theme_minimal() +
+   labs(title = "Product-wise Total Amount",
+       x = "Product", y = "Total Amount")
> # Plotting: Salesperson-wise total amount
> ggplot(salesperson_sales, aes(x = `Sales Person`, y = Total_Amount)) +
+   geom_bar(stat = "identity") +
+   theme_minimal() +
+   labs(title = "Salesperson-wise Total Amount",
+       x = "Sales Person", y = "Total Amount")
> # Plotting: Salesperson-wise average amount per box
> ggplot(salesperson_sales, aes(x = `Sales Person`, y = Average_Amount_per_Box)) +
+   geom_bar(stat = "identity") +
+   theme_minimal() +
+   labs(title = "Salesperson-wise Average Amount per Box",
+       x = "Sales Person", y = "Average Amount per Box")
>
> # Plotting: Distribution of Amount
> ggplot(data.df, aes(x = Amount)) +
+   geom_histogram(binwidth = 1000, fill = "blue", color = "black", alpha = 0.7) +
+   theme_minimal() +
+   labs(title = "Distribution of Amount",
+       x = "Amount", y = "Frequency")
>
```



Distribution of Amount

---

# Screenshot 2

**assignment-1 - master - RStudio**

File Edit Code View Plots Session Build Debug Profile Tools Help

sample_data_10mins

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

**Environment / History / Connections / Git / Tutorial**

Data
- country_sales — 6 obs. of 4 variables
- data.df — 1094 obs. of 7 variables
- dispersions — 1 obs. of 6 variables
- product_sales — 22 obs. of 4 variables
- salesperson_sales — 25 obs. of 4 variables
- sample_data_10mins — 1094 obs. of 6 variables
- tendencies — 1 obs. of 5 variables

**Console / Terminal / Background Jobs**

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```
+   geom_bar(stat = "identity") +
+   theme_minimal() +
+   labs(title = "Salesperson-wise Total Amount",
+       x = "Sales Person", y = "Total Amount")
> # Plotting: Salesperson-wise average amount per box
> ggplot(salesperson_sales, aes(x = `Sales Person`, y = Average_Amount_per_Box)) +
+   geom_bar(stat = "identity") +
+   theme_minimal() +
+   labs(title = "Salesperson-wise Average Amount per Box",
+       x = "Sales Person", y = "Average Amount per Box")
>
> # Plotting: Distribution of Amount
> ggplot(data.df, aes(x = Amount)) +
+   geom_histogram(binwidth = 1000, fill = "blue", color = "black", alpha = 0.7) +
+   theme_minimal() +
+   labs(title = "Distribution of Amount",
+       x = "Amount", y = "Frequency")
> # Plotting: Distribution of Amount per Box
> ggplot(data.df, aes(x = Amount_per_Box)) +
+   geom_histogram(binwidth = 10, fill = "green", color = "black", alpha = 0.7) +
+   theme_minimal() +
+   labs(title = "Distribution of Amount per Box",
+       x = "Amount per Box", y = "Frequency")
> |
```



Distribution of Amount per Box

**Screenshot 1 — RStudio (assignment-1 - master)**

Console code:

```
+      geom_bar(stat = "identity") +
+      theme_minimal() +
+      labs(title = "Salesperson-wise Average Amount per Box",
+           x = "Sales Person", y = "Average Amount per Box")
>
> # Plotting: Distribution of Amount
> ggplot(data.df, aes(x = Amount)) +
+      geom_histogram(binwidth = 1000, fill = "blue", color = "black", alpha = 0.7) +
+      theme_minimal() +
+      labs(title = "Distribution of Amount",
+           x = "Amount", y = "Frequency")
> # Plotting: Distribution of Amount per Box
> ggplot(data.df, aes(x = Amount_per_Box)) +
+      geom_histogram(binwidth = 10, fill = "green", color = "black", alpha = 0.7) +
+      theme_minimal() +
+      labs(title = "Distribution of Amount per Box",
+           x = "Amount per Box", y = "Frequency")
> # Scatter Plot: Amount vs. Boxes Shipped
> ggplot(data.df, aes(x = `Boxes Shipped`, y = Amount)) +
+      geom_point() +
+      theme_minimal() +
+      labs(title = "Amount vs. Boxes Shipped",
+           x = "Boxes Shipped", y = "Amount")
>
```

Plot title: Amount vs. Boxes Shipped (scatter plot of Amount vs Boxes Shipped)



**Screenshot 2 — RStudio (assignment-1 - master)**

Console code:

```
> ggplot(data.df, aes(x = Amount)) +
+      geom_histogram(binwidth = 1000, fill = "blue", color = "black", alpha = 0.7) +
+      theme_minimal() +
+      labs(title = "Distribution of Amount",
+           x = "Amount", y = "Frequency")
> # Plotting: Distribution of Amount per Box
> ggplot(data.df, aes(x = Amount_per_Box)) +
+      geom_histogram(binwidth = 10, fill = "green", color = "black", alpha = 0.7) +
+      theme_minimal() +
+      labs(title = "Distribution of Amount per Box",
+           x = "Amount per Box", y = "Frequency")
> # Scatter Plot: Amount vs. Boxes Shipped
> ggplot(data.df, aes(x = `Boxes Shipped`, y = Amount)) +
+      geom_point() +
+      theme_minimal() +
+      labs(title = "Amount vs. Boxes Shipped",
+           x = "Boxes Shipped", y = "Amount")
> # Box Plot: Amount per Box by Country
> ggplot(data.df, aes(x = Country, y = Amount_per_Box)) +
+      geom_boxplot() +
+      theme_minimal() +
+      labs(title = "Amount per Box by Country",
+           x = "Country", y = "Amount per Box")
>
```

Plot title: Amount per Box by Country (box plots for Australia, Canada, India, New Zealand, UK, USA)

File Edit Code View Plots Session Build Debug Profile Tools Help

sample_data_10mins

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

Environment History Connections Git Tutorial

R - Global Environment

Data
| country_sales | 6 obs. of 4 variables |
| data.df | 1094 obs. of 7 variables |
| dispersions | 1 obs. of 6 variables |
| product_sales | 22 obs. of 4 variables |
| salesperson_sales | 25 obs. of 4 variables |
| sample_data_10mins | 1094 obs. of 6 variables |
| tendencies | 1 obs. of 5 variables |

Console Terminal Background Jobs

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```r
> ggplot(data.df, aes(x = Amount_per_Box)) +
+     geom_histogram(binwidth = 10, fill = "green", color = "black", alpha = 0.7) +
+     theme_minimal() +
+     labs(title = "Distribution of Amount per Box",
+         x = "Amount per Box", y = "Frequency")
> # Scatter Plot: Amount vs. Boxes Shipped
> ggplot(data.df, aes(x = `Boxes Shipped`, y = Amount)) +
+     geom_point() +
+     theme_minimal() +
+     labs(title = "Amount vs. Boxes Shipped",
+         x = "Boxes Shipped", y = "Amount")
> # Box Plot: Amount per Box by Country
> ggplot(data.df, aes(x = Country, y = Amount_per_Box)) +
+     geom_boxplot() +
+     theme_minimal() +
+     labs(title = "Amount per Box by Country",
+         x = "Country", y = "Amount per Box")
> # Density Plot: Distribution of Amount per Box
> ggplot(data.df, aes(x = Amount_per_Box)) +
+     geom_density(fill = "blue", alpha = 0.7) +
+     theme_minimal() +
+     labs(title = "Density Plot of Amount per Box",
+         x = "Amount per Box", y = "Density")
> |
```

Files Plots Packages Help Viewer Presentation

Density Plot of Amount per Box

(Density plot: y-axis "Density" from 0.000 to 0.010, x-axis "Amount per Box" from 0 to 4000)

---

File Edit Code View Plots Session Build Debug Profile Tools Help

sample_data_10mins

| | Sales Person | Country | Product | Date | Amount | Boxes Shipped |
|---|---|---|---|---|---|---|
| 44 | Andria Kimpton | India | Spicy Special Slims | 2022-02-23 | 6307 | 142 |
| 45 | Madelene Upcott | Canada | Almond Choco | 2022-08-22 | 7602 | 102 |
| 46 | Kaine Padly | USA | Peanut Butter Cubes | 2022-02-16 | 6790 | 188 |
| 47 | Van Tuxwell | USA | 50% Dark Bites | 2022-01-13 | 9737 | 160 |
| 48 | Curtice Advani | Australia | Milk Bars | 2022-02-14 | 6979 | 18 |
| 49 | Roddy Speechley | India | Eclairs | 2022-06-10 | 4382 | 303 |
| 50 | Curtice Advani | India | Fruit & Nut Bars | 2022-07-07 | 5243 | 176 |
| 51 | Curtice Advani | Canada | Almond Choco | 2022-03-24 | 4865 | 70 |

Showing 44 to 52 of 1,094 entries, 6 total columns

Environment History Connections Git Tutorial

R - Global Environment

Data
| country_sales | 6 obs. of 4 variables |
| data.df | 1094 obs. of 7 variables |
| dispersions | 1 obs. of 6 variables |
| product_sales | 22 obs. of 4 variables |
| salesperson_sales | 25 obs. of 4 variables |
| sample_data_10mins | 1094 obs. of 6 variables |
| tendencies | 1 obs. of 5 variables |

Console Terminal Background Jobs

R 4.4.0 · C:/Users/Shaurya/Desktop/SEM 7 SUBJECTS/Big Data Visualization/Assignments/Assignment 1/Dataset/assignment-1/

```r
> ggplot(data.df, aes(x = `Boxes Shipped`, y = Amount)) +
+     geom_point() +
+     theme_minimal() +
+     labs(title = "Amount vs. Boxes Shipped",
+         x = "Boxes Shipped", y = "Amount")
> # Box Plot: Amount per Box by Country
> ggplot(data.df, aes(x = Country, y = Amount_per_Box)) +
+     geom_boxplot() +
+     theme_minimal() +
+     labs(title = "Amount per Box by Country",
+         x = "Country", y = "Amount per Box")
> # Density Plot: Distribution of Amount per Box
> ggplot(data.df, aes(x = Amount_per_Box)) +
+     geom_density(fill = "blue", alpha = 0.7) +
+     theme_minimal() +
+     labs(title = "Density Plot of Amount per Box",
+         x = "Amount per Box", y = "Density")
> # Plotting: Salesperson-wise total amount
> ggplot(salesperson_sales, aes(x = "", y = Total_Amount, fill = `Sales Person`)) +
+     geom_bar(stat = "identity", width = 1) +
+     coord_polar("y") +
+     theme_void() +
+     labs(title = "Salesperson-wise Total Amount")
> |
```

Files Plots Packages Help Viewer Presentation

Salesperson-wise Total Amount

Sales Person
- Andria Kimpton
- Barr Faughny
- Beverie Moffet
- Brien Boise
- Camilla Castle
- Ches Bonnell
- Curtice Advani
- Dennison Crosswaite
- Dotty Strutley
- Gigi Bohling
- Gunar Cockshoot
- Husein Augar
- Jan Morforth
- Jehu Rudeforth
- Kaine Padly
- Karlen McCaffrey
- Kelci Walkden
- Madelene Upcott
- Mallorie Waber
- Marney O'Breen
- Oby Sorrel
- Rafaelita Blaksland
- Roddy Speechley
- Van Tuxwell
- Wilone O'Kielt