CSE - 5334 - 001 : Data Mining

Assignment 1 - Exploratory Data Analysis - R

Instructor: Dr. Elizabeth D Diaz

Team 18: Urmi Manish Sheth – 1002064934

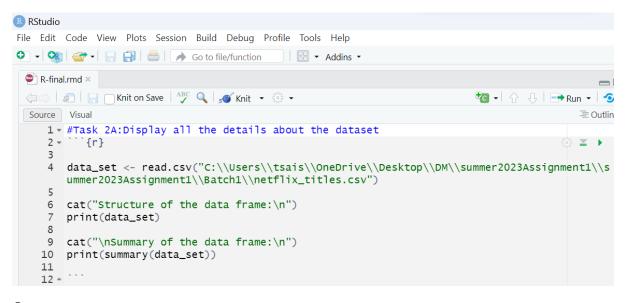
Harsh Navinbhai Shah – 1002057387

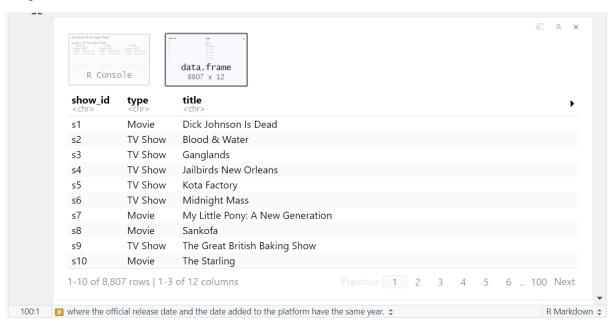
Sai Swetha Tadivaka – 1002112726

Part 2: R (40 Points) Perform all the tasks in R, using a different R notebook. The data and questions will be the same as those in Task 1.

Task 1A (2 points): Display all the details about the dataset.

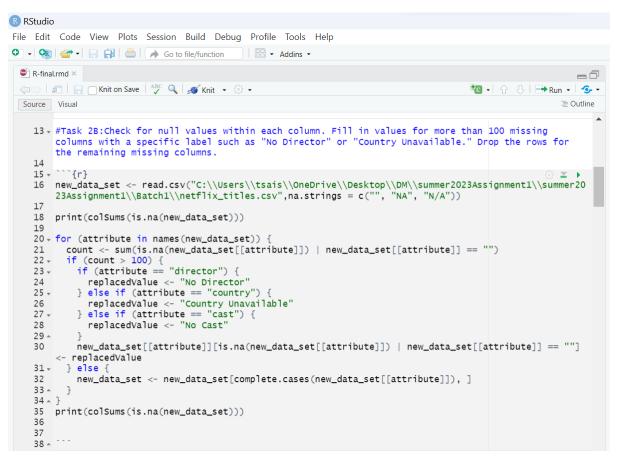
Solution





Task 2B (3 points): Check for null values within each column. Fill in values for more than 100 missing columns with a specific label such as "No Director" or "Country Unavailable." Drop the rows for the remaining missing columns.

Solution



Output: Before eliminating null values

```
show_id
                 type
                            title
                                      director
                                                      cast
       0
                   0
                               0
                                       2634
                                                      825
 country
           date_added release_year
                                       rating
                                                  duration
     831
                  10
listed_in
          description
```

Output: After eliminating null values

```
show_id
                   type
                               title
                                          director
                                                            cast
        0
                     0
                                   0
                                                 0
                                                               0
 country
            date_added release_year
                                            rating
                                                        duration
                     0
                                   0
                                                               0
        0
                                                 0
listed_in
           description
        0
```

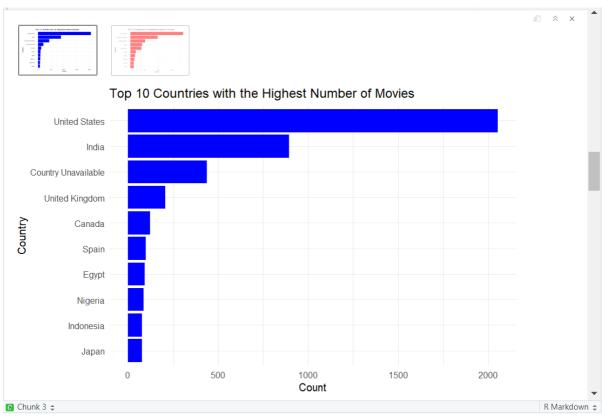
Task 2C (5 points): Create a horizontal bar chart displaying the top 10 countries with the highest number of movies and TV shows.

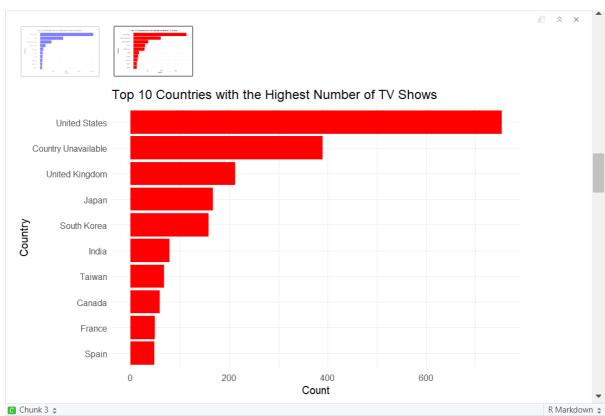
```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
R-final.rmd* ×

↓ □ | □ | □ Knit on Save | ♣ Q | ♣ Knit ▼ ② ▼
                                                                                                    🐿 - | 👉 🕛 | ➡ Run - | 🤣 -
 Source Visual
                                                                                                                       ■ Outline
     40 + ```{r}
     41
    42 # # Task 1C (5 points): Create a horizontal bar chart displaying the top 10 countries with
    43 # # the highest number of movies and TV shows.
    45 library(ggplot2)
    46 movies_data <- subset(new_data_set, type == "Movie")
     47 tv_shows_data <- subset(new_data_set, type == "TV Show")
    48
         top_movies_of_10_countries <- head(sort(table(movies_data$country), decreasing = TRUE), 10)</pre>
     49
         movies_data <- data.frame(</pre>
     50
            Country = names(top_movies_of_10_countries),
    51
           Count = as.numeric(top_movies_of_10_countries)
     52 )
     53 movies_data$Country <- factor(movies_data$Country, levels = rev(movies_data$Country))
    54 movies_horizonal_plot <- ggplot(movies_data, aes(x = Count, y = reorder(Country, Count))) +

55 geom_bar(stat = "identity", fill = "blue") +

1abs(x = "Count", y = "Country", title = "Top 10 Countries with the Highest Number of Movies") +
     57
            theme_minimal()
     58 top_tv_shows_countries <- head(sort(table(tv_shows_data$country), decreasing = TRUE), 10)
     59
         tv_shows_data <- data.frame(
    60
            Country = names(top_tv_shows_countries),
     61
            Count = as.numeric(top_tv_shows_countries)
     62 )
     63 tv_shows_data$Country <- factor(tv_shows_data$Country, levels = rev(tv_shows_data$Country))</p>
    tv_shows_horizontal_plot <- ggplot(tv_shows_data, aes(x = Count, y = reorder(Country, Count))) +
geom_bar(stat = "identity", fill = "red") +
labs(x = "Count", y = "Country", title = "Top 10 Countries with the Highest Number of TV Shows")</pre>
     67
           theme_minimal()
    68 par(mfrow = c(1, 2))
     69
         print(movies_horizonal_plot)
     70 print(tv_shows_horizontal_plot)
```



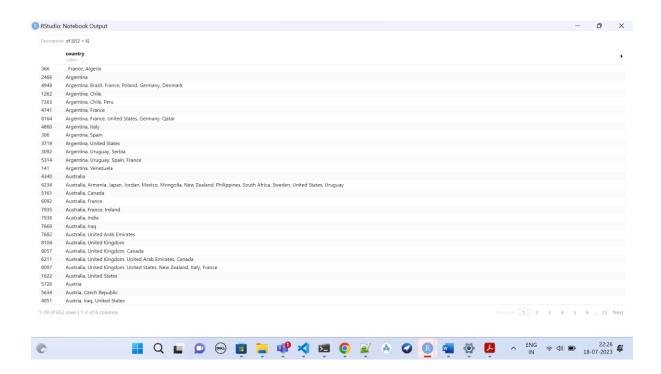


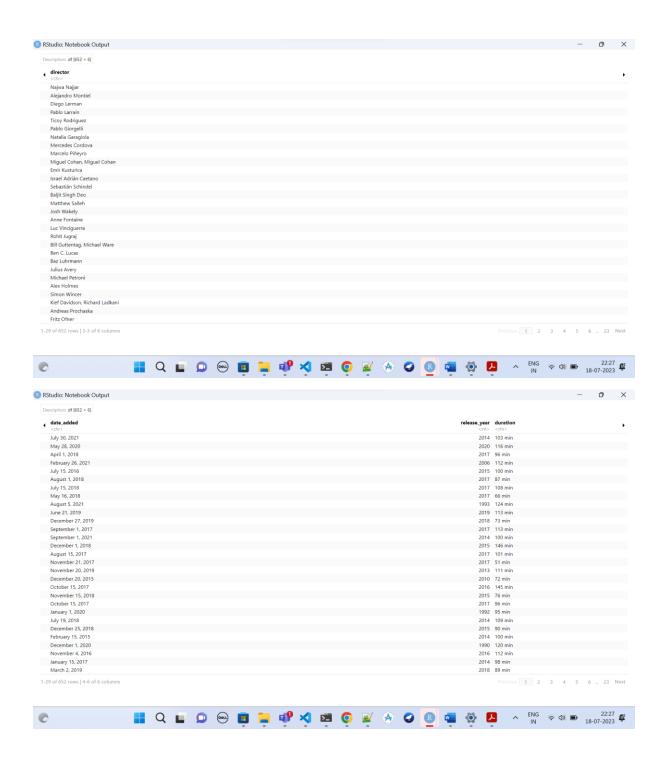
Task 2D (4 points): Print the first row based on the longest duration time of a movie from each country. Include information such as the director, date added, release year, duration, and description of the movie.

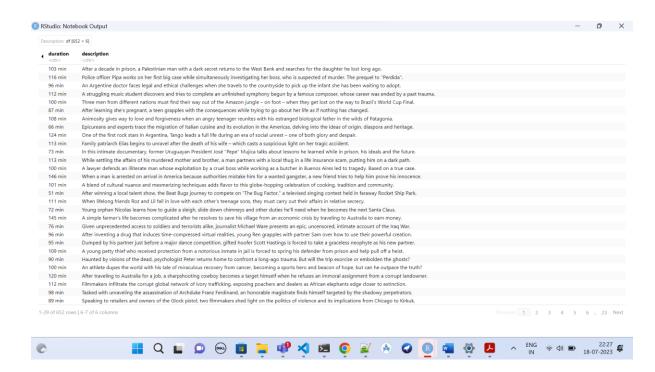
Solution

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help

⟨□□⟩ | Æ□ | ☐ Knit on Save | △ABC | ✓ Knit ▼ ∅ ▼
                                                                                      🐿 - | û 🖖 | ➡ Run - | 🤣 -
 Source Visual
                                                                                                       ■ Outline
    73 + ```{r}
    74 # Task 1D (4 points): Print the first row based on the longest duration time of a movie from each
        country. Include information such as the director, date added, release year, duration, and
        description of the movie.
    76 movies_data <- subset(new_data_set, type == "Movie")
77 movies_data$duration_in_mins <- as.integer(sub(" min", "", movies_data$duration))
78 vidx <- unlist(by(movies_data, movies_data$country, function(requiredData) {
        index <- which.max(requiredData$duration_in_mins)</pre>
    79
    80
          rownames(requiredData)[index]
    81 4 }))
    85 🛎
```







Task 2E (4 points): Display the title, director, date added, and release date of movies where the official release date and the date added to the platform have the same year.

Solution

```
Refinal.rmd* X

Source Visual

86 * # Task 1E (4 points): Display the title, director, date added, and release date of movies
87 * # where the official release date and the date added to the platform have the same year.

88

89 * ``{r}

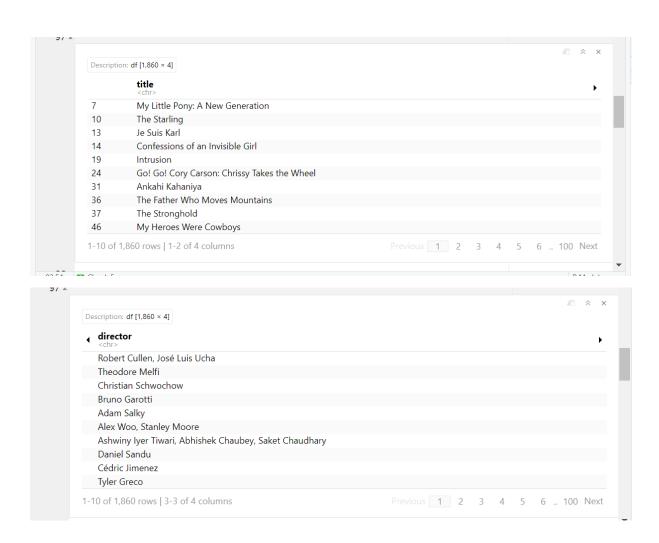
90

91 movies_data_set <- subset(new_data_set, type == "Movie")
92 movies_data_set$date_added <- as.Date(movies$date_added, format = "%B %d, %Y")
93 filtered_movies <- subset(movies_data_set, format(date_added, "%Y") == as.character(release_year))
94 print(filtered_movies[, c('title', 'director', 'date_added', 'release_year')])

95

96

97 * ```
```



Description: df [1,860 × 4]		
•	date_added <date></date>	release_ye <ii< th=""></ii<>
	2021-09-24	20
	2021-09-24	20
	2021-09-23	20
	2021-09-22	20
	2021-09-22	20
	2021-09-21	20
	2021-09-17	20
	2021-09-17	20
	2021-09-17	20
	2021-09-16	20

Task 2F (4 points): Display the director, release year, and the number of movies and TV shows directed by each director within a year. Sort the results in descending order based on count

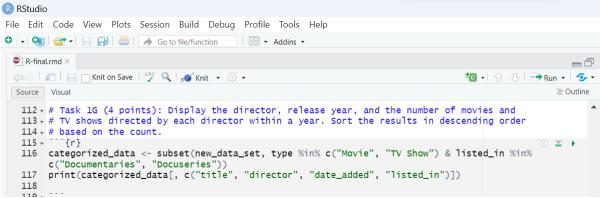
Solution

```
99 Task 1F (4 points): Display the director, release year, and the number of movies and
101 # Task 1F (4 points): Display the director, release year, and the number of movies and
102 # TV shows directed by each director within a year. Sort the results in descending order
103 # based on the count.
104
105 movies_tvshows_data <- subset(new_data_set, type %in% c("Movie", "TV Show"))
106 director_count <- aggregate(cbind(count = type) ~ director + release_year, data = movies_tvshows_data, FUN = length)
107 director_count <- director_count[order(-director_count$count), ]
108 desired_columns <- c("director", "release_year", "count")
109 print(director_count[desired_columns])
110
```

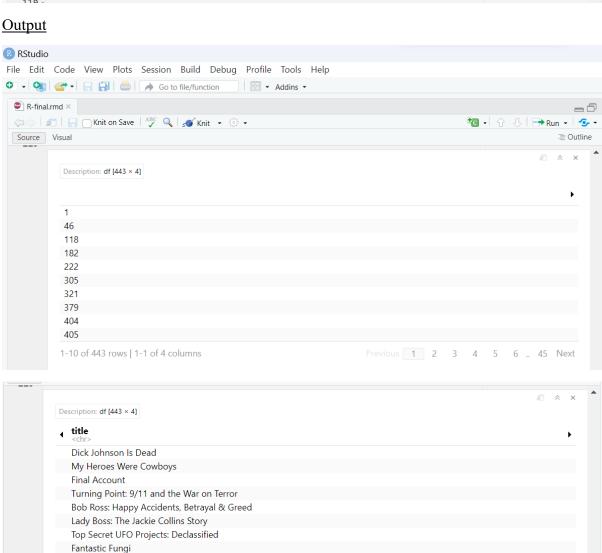


Task 2G (3 points): Display the title, director, date added, and category of movies or TV shows from the Documentary/Docuseries category.

Solution

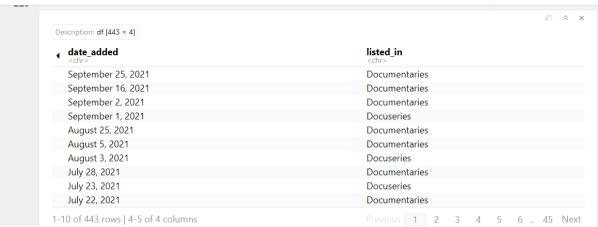


The Movies That Made Us 9to5: The Story of a Movement 1-10 of 443 rows | 2-2 of 4 columns

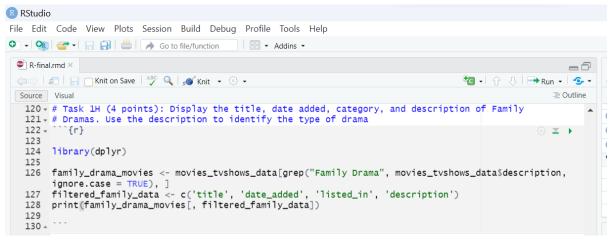


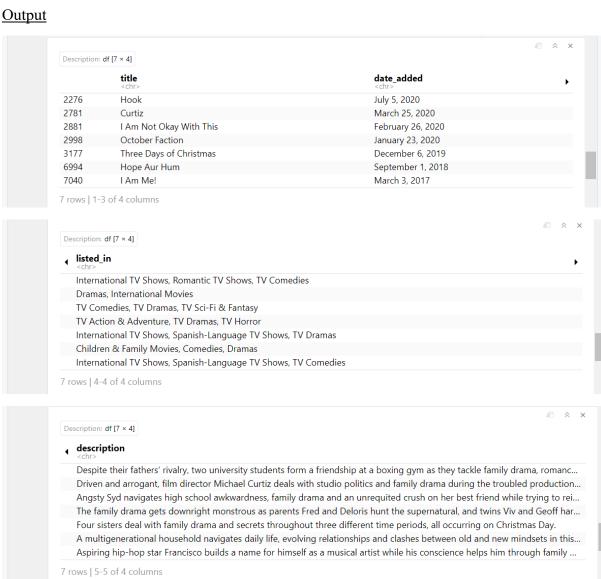
Previous 1 2 3 4 5 6 ... 45 Next





Task 2H (4 points): Display the title, date added, category, and description of Family Dramas. Use the description to identify the type of drama.





Task 2I (5 points): Plot the distribution of TV shows based on the number of seasons using a horizontal bar chart. Group the seasons into the following categories:

Less than 3 seasons

3 seasons

4 seasons

5 to less than 10 seasons

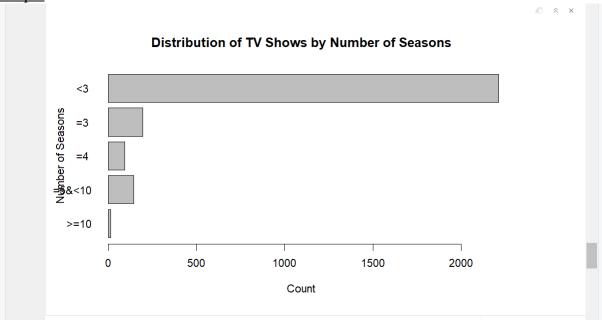
10 or more seasons

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
• Go to file/function
                                                       □ - Addins -
  R-final.rmd ×
  🕶 - | ↑ 👵 | ➡ Run - | 🤣 -
  Source Visual

    Outline

    131 - # Task 1I (5 points): Plot the distribution of TV shows based on the number of seasons 132 - # using a horizontal bar chart. Group the seasons into the following categories:
    133 - # Less than 3 seasons
    134 + # 3 seasons
    135 - # 4 seasons
    136 - # 5 to less than 10 seasons
    137 - # 10 or more seasons
    138
    139 - ```{r}
    140
    141 tv_shows <- subset(new_data_set, type == "TV Show")
142 season_counts <- table(gsub("Seasons?", "", tv_shows$duration))
    143 categories <- c(
            '<3 ' = sum(season_counts[as.numeric(names(season_counts)) < 3]),
'=3 ' = sum(season_counts[as.numeric(names(season_counts)) == 3]),
'=4 ' = sum(season_counts[as.numeric(names(season_counts)) == 4]),
    144
    145
    146
           '>=5&<10' = sum(season_counts[as.numeric(names(season_counts)) >= 5 &
    147
         as.numeric(names(season_counts)) < 10])
    148
              '>=10' = sum(season_counts[as.numeric(names(season_counts)) >= 10])
    149
    barplot(rev(categories), horiz = TRUE, xlab = 'Count', ylab = 'Number of Seasons',

main = 'Distribution of TV Shows by Number of Seasons', las = 1)
    152
    153
    154 -
    155
```



Task 2J (6 points): Display a side-by-side pie chart showing the distribution of movie and TV show ratings.

Movie Ratings:

Uncut/Not rated

Restricted

Parental guidance

General audience

Adults only

TV Show Ratings:

All Children

Older Children

Parental Presence

General audience

Mature

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
*C • | ↑ ↓ | → Run • | * •

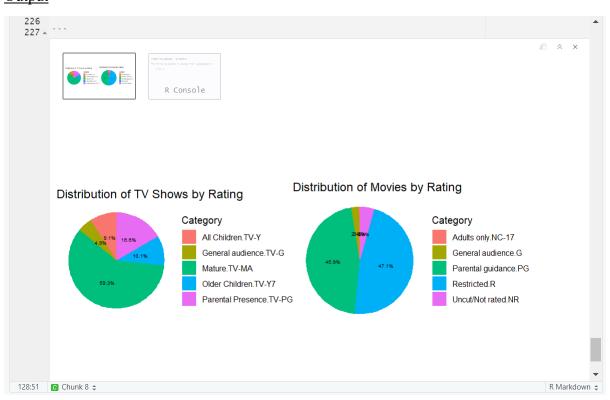
⟨□□⟩ | Æ□ | ☐ Knit on Save | ♣ Q | ★ Knit ▼ ♠ ▼
   Source Visual

    Outline

     158 # Task 1J (6 points): Display a side-by-side pie chart showing the distribution of movie
      160 # Movie Ratings:
     161 # Uncut/Not rated
162 # Restricted
163 # Parental guidance
164 # General audience
      165 # Adults only
     166 # TV Show Ratings:
167 # All Children
     168 # Older Children
169 # Parental Presence
     170 # General audience
     171 # Mature
172 library(ggplot2)
173 library(gridExtra)
     174
     175 tv_shows <- subset(new_data_set, type == "TV Show")
176 tv_shows_rating_count <- table(tv_shows$rating)
177 movies <- subset(new_data_set, type == "Movie")
178 movies_rating_count <- table(movies$rating)
     179 tv_shows_categories <- c(
                 'Shows_tategories <= C(
'All Children' = tv_shows_rating_count['TV-Y'],
'Older Children' = tv_shows_rating_count['TV-Y7'] + tv_shows_rating_count['TV-Y7-FV'],
'Parental Presence' = tv_shows_rating_count['TV-PG'],
'General audience' = tv_shows_rating_count['TV-G'],
'Mature' = tv_shows_rating_count['TV-MA']</pre>
     180
     181
     182
     183
     184
     185 )
```

```
186 movies_categories <- c(
         'Uncut/Not rated' = movies_rating_count['NR'],
187
         'Restricted' = movies_rating_count['R'],
188
        'Parental guidance' = movies_rating_count['PG'] + movies_rating_count['PG-13'],
'General audience' = movies_rating_count['G'],
189
190
191
         'Adults only' = movies_rating_count['NC-17']
192 )
193
194 tv_shows_data <- data.frame(
        Category = names(tv_shows_categories),
Count = as.numeric(tv_shows_categories)
195
196
197
198
    movies_data <- data.frame(
199
        Category = names(movies_categories);
200
        Count = as.numeric(movies_categories)
201
202
```

```
202
     tv_shows_data$Percentage <- tv_shows_data$Count / sum(tv_shows_data$Count) * 100 movies_data$Percentage <- movies_data$Count / sum(movies_data$Count) * 100
203
204
205
     tv_shows_plot <- ggplot(tv_shows_data, aes(x = "", y = Count, fill = Category)) +
    geom_bar(stat = "identity", width = 1) +
    coord_polar("y", start = 0) +
    labs(title = "Distribution of TV Shows by Rating") +</pre>
206
207
208
209
210
           theme_void() +
           geom_text(aes(label = paste0(sprintf("%.1f", Percentage), "%")),
211
                          position = position_stack(vjust = 0.5),
size = 2) # Adjust the size here
212
213
214
movies_horizontal_plot <- ggplot(movies_data, aes(x = "", y = Count, fill = Category)) +
geom_bar(stat = "identity", width = 1) +
coord_polar("y", start = 0) +
labs(title = "Distribution of Movies by Rating") +
219
          theme_void() +
220
          geom_text(aes(label = paste0(sprintf("%.1f", Percentage), "%")),
221
                          position = position_stack(vjust = 0.5),
222
                          size = 2) # Adjust the size here
224 grid.arrange(tv_shows_plot, movies_horizontal_plot, ncol = 2)
225
226
227 🛦
```



References:

R Language Definition (r-project.org)