



Exploring Movies

April 27, 2022

The dataset used for this challenge is titled 'movies.csv'.

What topic does the dataset cover?

This dataset contains information about the top 1000 highest grossing holywood films. It is up to date as of 10th January 2022.

Acknowledgements

This data has been scraped from multiple site and has been added together for performing various datat operations. The data has been taken from idmb, rotten tomatoes and many other sites. It can be found at <https://www.kaggle.com/datasets/sanjeetsinghnaik/top-1000-highest-grossing-movies>.

Variables

- title: Movie name
- movie_info: Information on the given movie
- distributor: distributor for movie
- release_date: When the movie was released
- domestic_sales: Domestic sales in dollars
- international_sales: International sales in dollars
- world_sales: World sales in dollars
- runtime: runtime of movie in minutes
- licence: licence of movie (G, PG, etc.)
- genre: Main genre of movie

Assignment

Complete each of the following questions using R Studio and submit your answers as a detailed report.

1. Visualize and describe the world sales of movies in millions of dollars.
2. Visualize and describe the distribution for movie runtimes in minutes.
3. Summarize how many movies each distributor has made.
 - (a) Which 3 distributors have made the most movies.
 - (b) Which 3 distributors have made the least movies?

4. Which three distributors have the highest average world sales?
5. Look at the world sales across the 5 distributors with the most amount of movies made. Does there appear to be a significant difference?
6. Count the number of movies that are in each genre. What is problematic about the way the data is coded?
7. Compare movie runtimes across the 4 genres that have made the most movies. Does there appear to be a significant difference? Use pairwise boxplots to support your answer.
8. Compare domestic sales across the 4 genres that produce the most movies. Does there appear to be a significant difference? Support your answer using pairwise boxplots.
9. Is there a linear relationship between the runtime of a movie and its world sales? Support your answer using a scatterplot.
10. Is there a linear relationship between the international sales of a movie and the domestic sales of a movie? Support your answer using a scatterplot.
11. If a movie made \$750 million domestically, what would you predict it to make internationally?
12. Count the number of movies in each license category. Is there anything problematic with the data?
13. Summarize the runtime of movies across the various licence categories. Does there appear to be any significant difference?
14. What is your favourite movie and why?

Solution:

```
1  ### Load packages and dataset
2  library(tidyverse)
3  library(ggplot2)
4  mov <- read.csv('/data/datasets/movies.csv')
5
6  ### Take a look at the dataset
7  glimpse(mov)
8
9  ### Question 1: Visualize the world sales of movies in millions of dollars
10 mov %>% mutate(mill = world_sales/1000000) %>%
11   ggplot(aes(x=mill)) +
12   geom_histogram(col='black', fill='green') +
13   labs(x="International Sales (millions of dollars)") +
14   theme_classic()
15
16 ### Question 2: Visualize the distribution for movie runtime in hours
17 mov %>% ggplot(aes(x=runtime_minutes)) +
18   geom_histogram(col='black', fill='red') +
19   labs(x="Runtime in minutes") +
20   theme_classic()
21
22 ### Question 3A: Which 3 distributors make the most movies?
23 mov %>% group_by(distributor) %>%
24   summarise(count=n()) %>%
25   arrange(desc(count)) %>% head(10)
26
27 ### Question 3B: Which 3 distributors make the least movies?
28 mov %>% group_by(distributor) %>%
```

```

29 summarise(count=n()) %>%
30 arrange(count) %>% head(10)
31
32 ### Question 4: Which three distributors have the highest average world sales?
33 mov %>% group_by(distributor) %>%
34 summarise(mean = mean(world_sales)) %>%
35 arrange(desc(mean))
36
37 ### Question 5: International sales across 5 biggest distributors
38 top <- c("Warner Bros.", "Walt Disney Studios Motion Pictures",
39          "Universal Pictures", "Twentieth Century Fox",
40          "Sony Pictures Entertainment (SPE)")
41
42 mov %>% filter(distributor %in% top) %>%
43 ggplot(aes(y=world_sales, fill=distributor)) +
44 geom_boxplot() +
45 theme_classic()
46
47 ### Question 6: Number of movies in each genre
48 mov %>% group_by(genre) %>%
49 summarise(count=n()) %>%
50 arrange(desc(count))
51
52 ### Question 7: Compare runtimes across top 4 genres
53 top <- c("Action", "Adventure", "Comedy", "Drama")
54 mov %>% filter(genre %in% top) %>%
55 ggplot(aes(y=runtime_minutes, fill=genre)) +
56 geom_boxplot() +
57 theme_classic()
58
59 ### Question 8: Compare domestic sales across top 4 genres
60 top <- c("Action", "Adventure", "Comedy", "Drama")
61 mov %>% filter(genre %in% top) %>%
62 ggplot(aes(y=domestic_sales, fill=genre)) +
63 geom_boxplot() +
64 theme_classic()
65
66 ### Question 9: Is there a relationship between runtime and world sales?
67 mov %>% ggplot(aes(x=runtime_minutes, y=world_sales)) +
68 geom_point() +
69 theme_classic()
70
71 ### Question 10: Is there a relationship between international sales and domestic
72 sales?
73 mov %>% ggplot(aes(x=domestic_sales, y=international_sales)) +
74 geom_point() +
75 geom_smooth() +
76 theme_classic()
77
78 ### Question 11: Eyeball prediction using smoothed line
79
80 ### Question 12: Count number of movies in each license category
81 mov %>% group_by(license) %>%
82 summarise(count = n())
83
84 ### Question 13: compare the runtime for movies in each license category
85 mov %>% ggplot(aes(y=runtime_minutes, fill=license))+
86 geom_boxplot() +
87 theme_classic()
88

```

