

Exercise:

- Analyse dataset movies
- Visualize main features of this dataset using ggplot2 package
- Create a chart with a few panels characterising 3 most important features of this dataset.

I. Data overview

title	year	length	budget
Length:58788	Min. :1893	Min. : 1.00	Min. : 0
Class :character	1st Qu.:1958	1st Qu.: 74.00	1st Qu.: 250000
Mode :character	Median :1983	Median : 90.00	Median : 3000000
	Mean :1976	Mean : 82.34	Mean : 13412513
	3rd Qu.:1997	3rd Qu.: 100.00	3rd Qu.: 15000000
	Max. :2005	Max. :5220.00	Max. :200000000
			NA's :53573

rating	votes	r1	r2
Min. : 1.000	Min. : 5.0	Min. : 0.000	Min. : 0.000
1st Qu.: 5.000	1st Qu.: 11.0	1st Qu.: 0.000	1st Qu.: 0.000
Median : 6.100	Median : 30.0	Median : 4.500	Median : 4.500
Mean : 5.933	Mean : 632.1	Mean : 7.014	Mean : 4.022
3rd Qu.: 7.000	3rd Qu.: 112.0	3rd Qu.: 4.500	3rd Qu.: 4.500
Max. :10.000	Max. :157608.0	Max. :100.000	Max. :84.500

r3	r4	r5	r6
Min. : 0.000	Min. : 0.000	Min. : 0.000	Min. : 0.00
1st Qu.: 0.000	1st Qu.: 0.000	1st Qu.: 4.500	1st Qu.: 4.50
Median : 4.500	Median : 4.500	Median : 4.500	Median :14.50
Mean : 4.721	Mean : 6.375	Mean : 9.797	Mean :13.04
3rd Qu.: 4.500	3rd Qu.: 4.500	3rd Qu.: 14.500	3rd Qu.:14.50
Max. :84.500	Max. :100.000	Max. :100.000	Max. :84.50

r7	r8	r9	r10
Min. : 0.00	Min. : 0.00	Min. : 0.000	Min. : 0.00
1st Qu.: 4.50	1st Qu.: 4.50	1st Qu.: 4.500	1st Qu.: 4.50
Median : 14.50	Median : 14.50	Median : 4.500	Median : 14.50
Mean : 15.55	Mean : 13.88	Mean : 8.954	Mean : 16.85
3rd Qu.: 24.50	3rd Qu.: 24.50	3rd Qu.: 14.500	3rd Qu.: 24.50
Max. :100.00	Max. :100.00	Max. :100.000	Max. :100.00

mpaa	Action	Animation	Comedy
Length:58788	Min. :0.00000	Min. :0.00000	Min. :0.00000
Class :character	1st Qu.:0.00000	1st Qu.:0.00000	1st Qu.:0.00000
Mode :character	Median :0.00000	Median :0.00000	Median :0.00000
	Mean :0.07974	Mean :0.06277	Mean :0.2938
	3rd Qu.:0.00000	3rd Qu.:0.00000	3rd Qu.:1.00000
	Max. :1.00000	Max. :1.00000	Max. :1.00000

Drama	Documentary	Romance	Short
Min. :0.000	Min. :0.00000	Min. :0.0000	Min. :0.0000
1st Qu.:0.000	1st Qu.:0.00000	1st Qu.:0.0000	1st Qu.:0.0000
Median :0.000	Median :0.00000	Median :0.0000	Median :0.0000
Mean :0.371	Mean :0.05906	Mean :0.0807	Mean :0.1609
3rd Qu.:1.000	3rd Qu.:0.00000	3rd Qu.:0.0000	3rd Qu.:0.0000
Max. :1.000	Max. :1.00000	Max. :1.0000	Max. :1.0000

The internet movie database imdb.com is a website devoted to collecting movie data supplied by studios and fans. The data frame contains 28819 rows and 24 variables:

- title: Title of the movie
- year: Year of release
- budget: Total budget (if known) in US dollars
- length: Duration of the movies in minutes
- rate: Average IMDB user rating
- votes: Number of IMDB user who rated this movie
- r1-10: Multiplying by ten gives percentile (to nearest 10%) of users who rated this movie a 1.
- mpaa: MPAA rating
- action, animation, comedy, drama, documentary, romance, short: Binary variables representing if a movie was classified as belonging to that genre.

Notably, the “budget” column contains missing data as there are movies’ budgets that are unknown. In addition, “mpaa” column also has a lot of blank cells.

II. Data processing

1. Create “category” column
2. Drop binary category columns
3. Edit data type of the remaining variables
4. Create “unknown” value for missing value in the “budget” column

After these steps, we have a data frame as follow:

5. Create sub table for rating and votes of each movie (rate_votes)

6. Create sub table for number and rating of category (genre_summary)

category	count	avg_rate	avg_votes
Action	4688	5.292022	2086.72440
Animation	3606	6.590627	311.25901
Comedy	14269	5.872612	730.96979
Documentary	3183	6.656645	85.81464
Drama	16952	6.189352	689.87394
Other	12786	5.448045	285.18622
Romance	580	6.058276	264.96379
Short	2724	6.287372	21.49339

7. Create sub table about rates and number of movies each year (yearly_trends)

year	avg_rate	movie_count
1893	7.000000	1
1894	4.888889	9
1895	5.500000	3
1896	5.269231	13
1897	4.677778	9
1898	5.040000	5

8. Create sub table about movies duration and movies rating (length_analysis)

- Short: Duration under 60 mins
- Medium: Duration between 60 mins and 120 mins

- Long: Duration more than 120 mins

length_category	avg_rate	movie_count
Long	6.792058	3752
Medium	5.738166	44220
Short	6.430741	10816

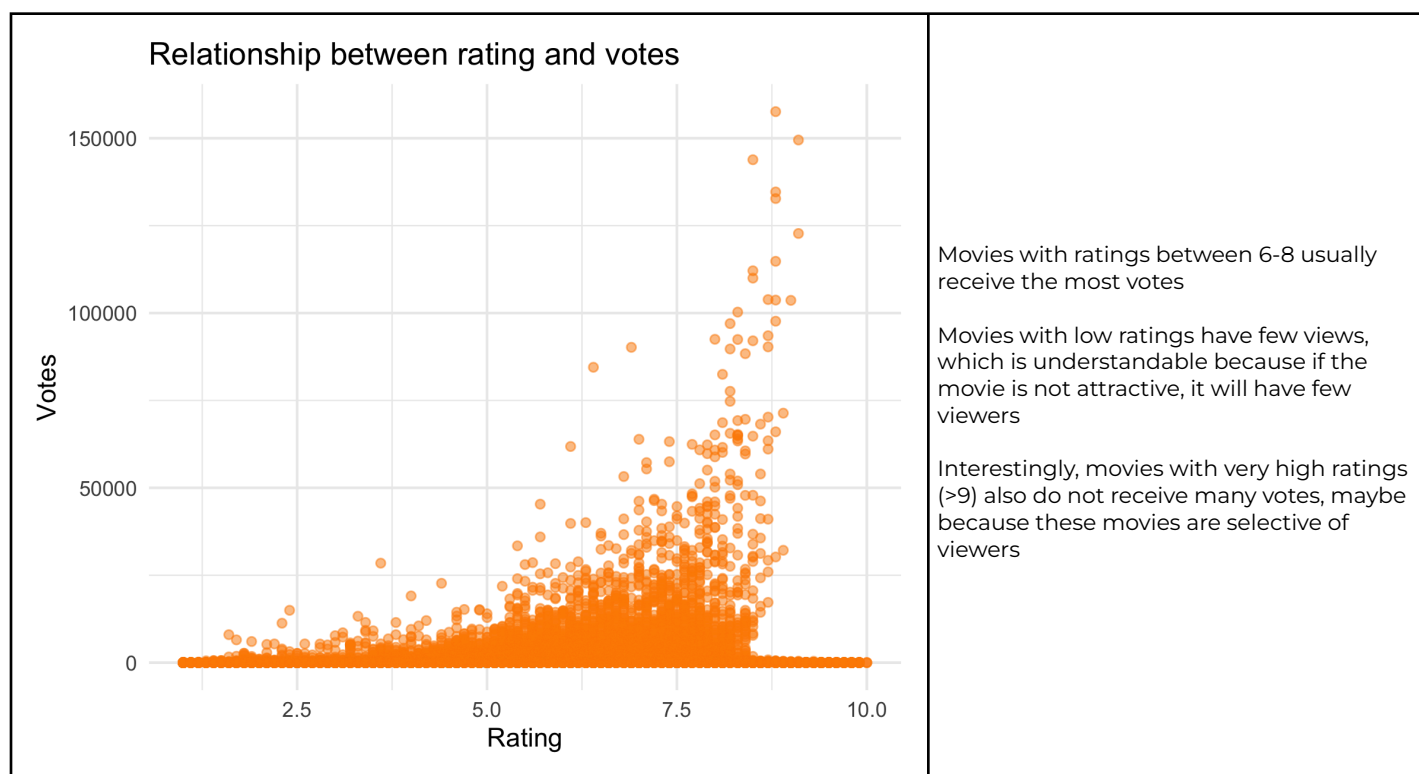
9. Create sub table for top movies in the dataset (top_movies)

title	year	rating	votes
Dimensia Minds Trilogy: The Hope Factor	2004	10.0	5
Fishing for Love	2001	10.0	5
Summer Sonata, A	2004	10.0	5
Drifting	2004	9.9	41
Himala	1982	9.9	26
Ivan Groznyy III	1988	9.9	26
New World, The	1982	9.9	24
Dimensia Minds Trilogy: The Reds	2004	9.9	15
Titanic vals	1964	9.9	13
Plight of Clownana, The	2004	9.9	11

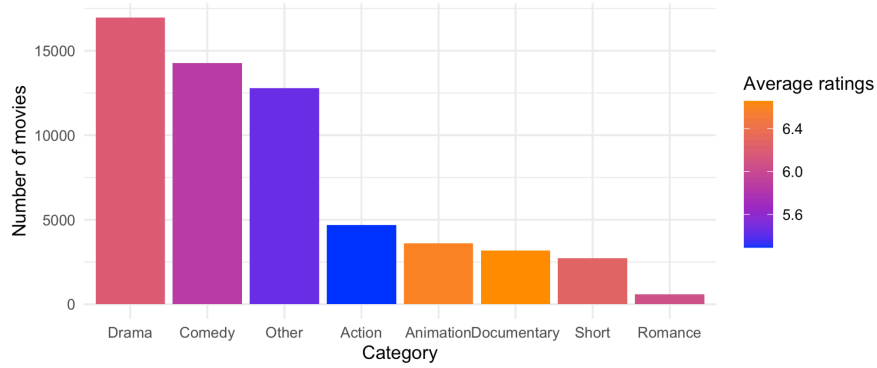
10. Create sub table for top movies in the dataset with votes more than 30 (top_movies_filtered)

title	year	rating	votes
Drifting	2004	9.9	41
Cashback	2004	9.8	51
Glissando	1985	9.8	37
Goodnite Charlie	2005	9.8	34
Tis kakomoiras	1963	9.6	86
Classic Queen	1992	9.6	39
Da no tien gu	1965	9.5	94
Punto y raya	2004	9.5	62

III. Data analysis



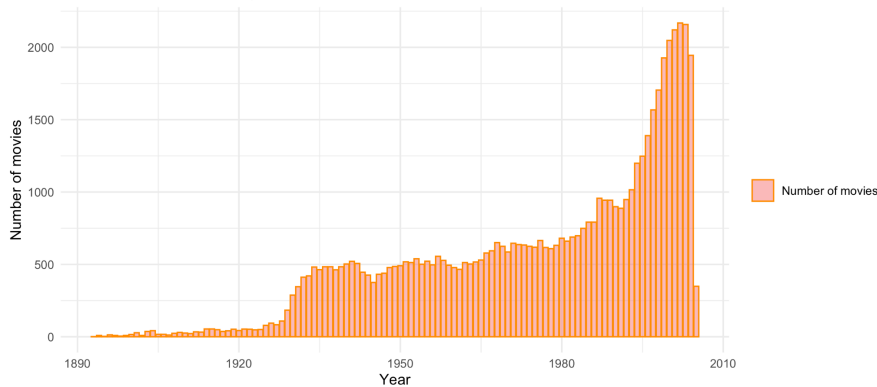
Number of movies and average rating of each category



It can be seen that drama and comedy are the two most popular genres, receiving the most ratings

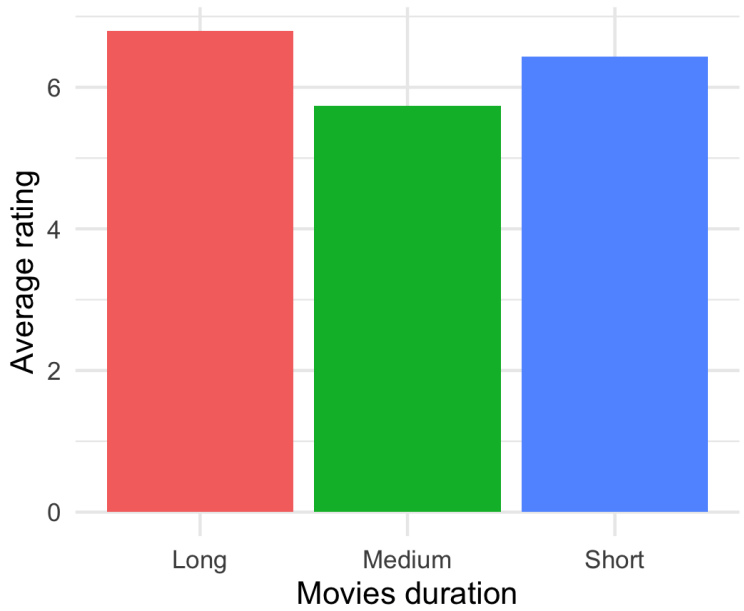
The genres that receive the highest ratings are animation and documentary, short

Number of movies over years



The number of films increased steadily over the years and increased significantly after the 2000s, reflecting the development of the film industry during this period.

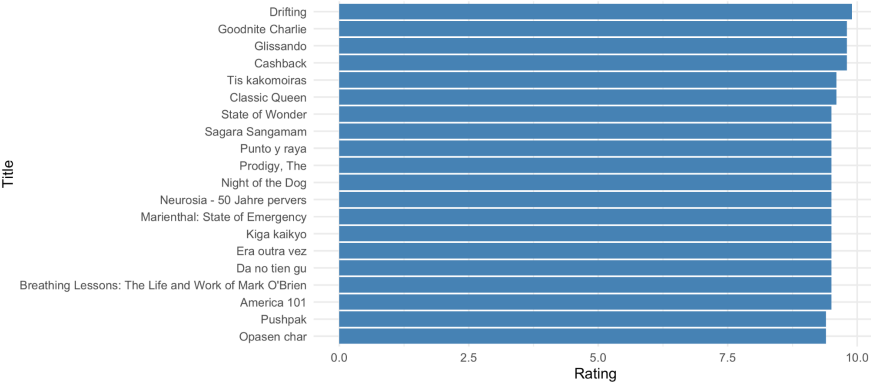
Average rating by movies duration



Medium length films (60-120 minutes) are generally rated highest.

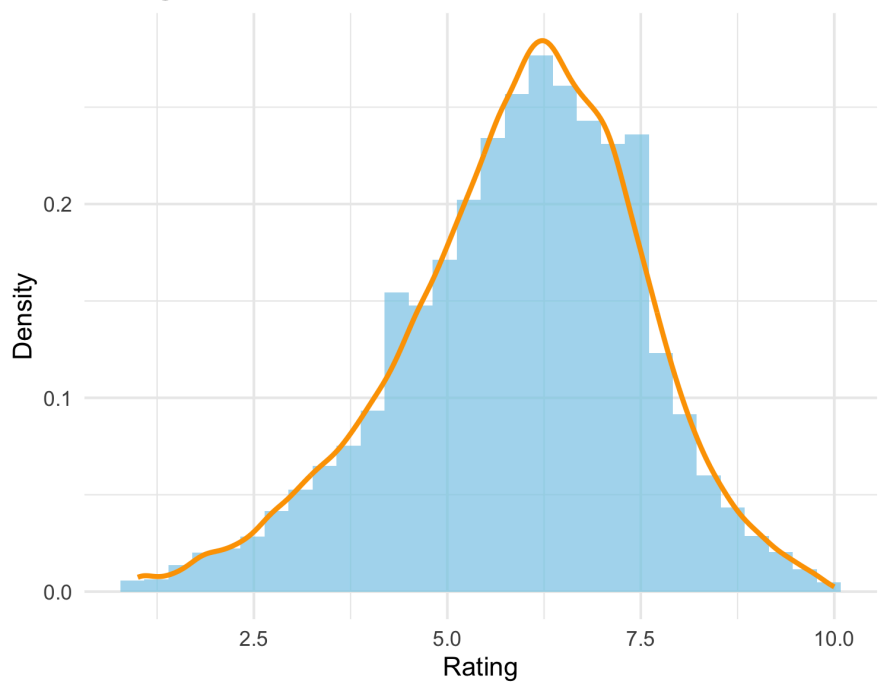
Short films (<60 minutes) are rated lower, possibly due to less engaging content.

Top 20 movies with highest rating (votes >= 30)



The chart shows the top 20 highest rated movies (with at least 30 votes)

Rating distribution



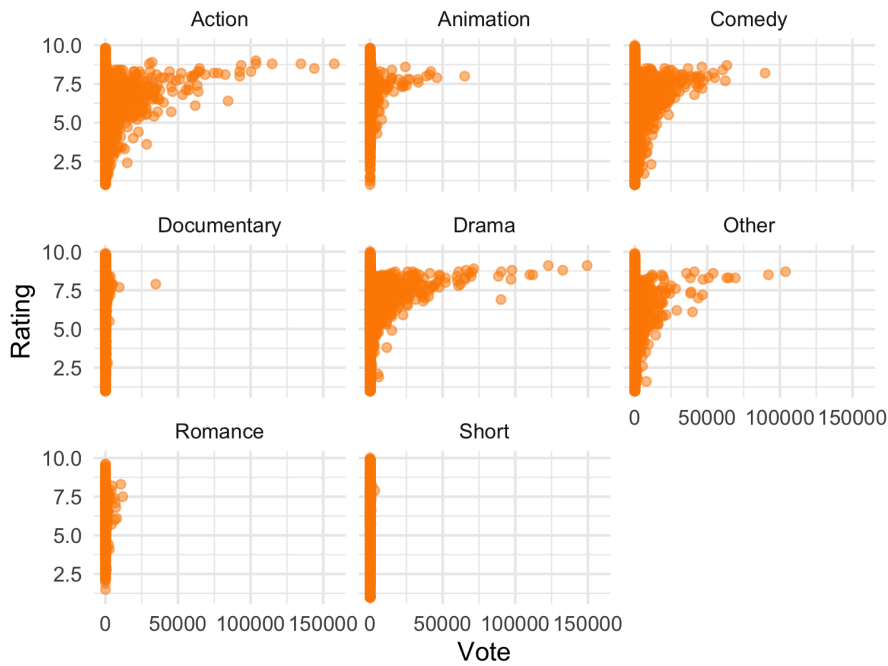
It can be seen that the most popular movies are rated at 6-7.

Correlation between numeric variables



It can be seen that there is no significant negative or positive correlation between the numerical variables.

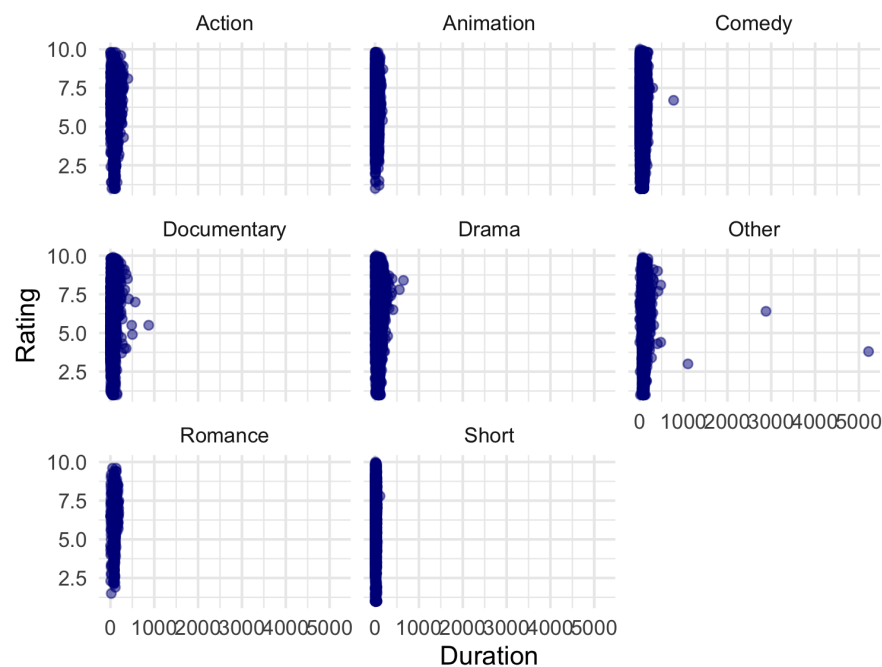
Relationship between votes and rates by category



As can be seen, the genres with the strongest correlation between rating and votes are action, drama, followed by comedy and unspecified genres respectively.

It can be seen that short films often receive very few votes.

Relationship between duration and rates by category



There isn't really a special correlation between film length and film rating. It's understandable since films are often rated on content, not length.