**Bern University of Applied Sciences**

Business Information Systems

# DSWI

## Animation and interaction

## Midterm Group Work



| | |
|---|---|
| Study programme | **BSc Business Information Systems** |
| Submitted by | **Birnbaumer Jamie-Lee, Djerrah Yacine, Zimmermann Yannik** |
| Subject area | **Data Science** |
| Expert | **Prof. Dr. Léonard Kwuida** |
| Datum of submission | **2. May 2021** |

# Content

# 1. Preparations

```
1  ## Loading librarys ##
2  library(readr)
3  library(tidyverse)
4  library(gifski)
5  library(ggplot2)
6  library(gganimate)
7  library(dplyr)
8  #Only used for masking count
9  library(plyr)
```

First, the necessary libraries must be initialised. Missing libraries can be installed via the Package Installer. To reproduce our tutorial, you will need the libraries listed on the left. Afterwards we will import our datasets. The source of the data is listed in the source directory. Generally, we work with a dataset from the movie database IMDb. More precisely, we use the data collections with all movie titles over more than 100 years as well as the corresponding movie ratings. We will be Inner-Joining the ImdbTitleRatings into ImdbTitleBasics Chose Inner-Join to make sure you will only insert datasets with an actual ratings value and skip the empty entries.

```
11 ## Importing Data ##
12 ImdbTitleBasics <- read_csv("ImdbTitleBasics.csv")
13 ImdbTitleRatings <- read_csv("ImdbTitleRatings.csv")
14
15 ## Joining Datasets ##
16 ImdbTitleBasicsAndRatings = inner_join(ImdbTitleBasics, ImdbTitleRatings, by = c("tconst" = "tconst"))
```

The following code excerpt shows how the data is prepared to be displayed in the plots. To maintain a significance, the data is reduced to the 10 most common genres & types.
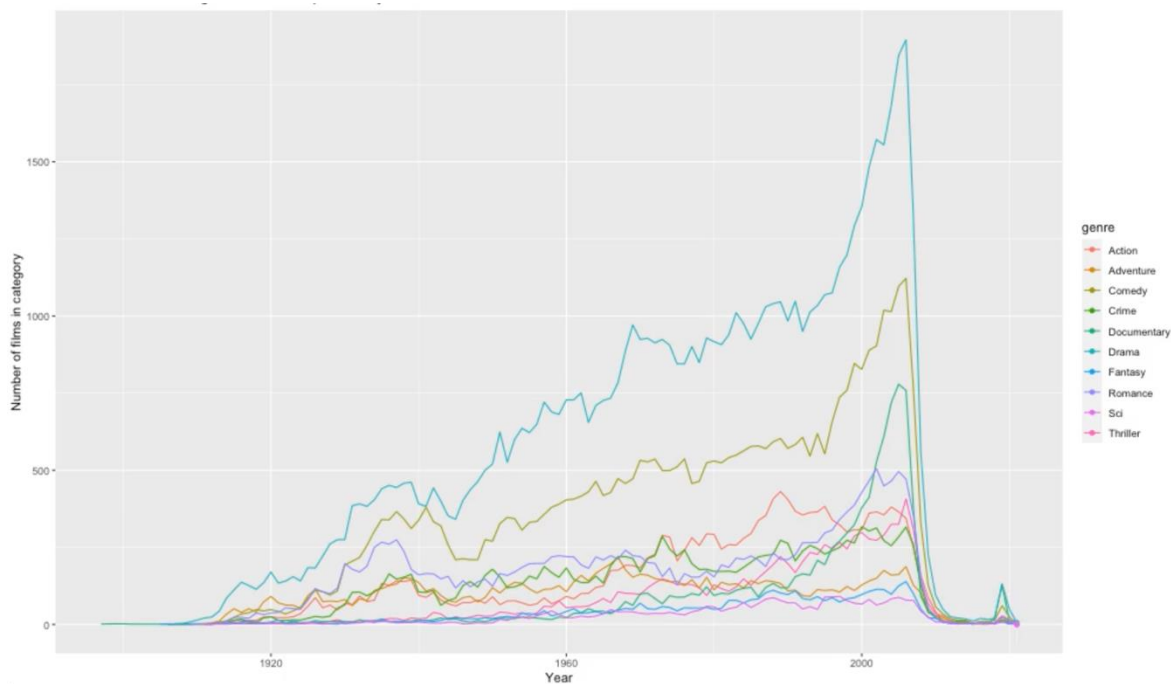
```
18 ## Preparing Data for plotting
19 # Filters for the dataset #
20 genre_filter = c("Documentary", "Drama", "Fantasy", "Action", "Romance", "Comedy", "Adventure", "Crime", "Sci", "Thriller")
21 data_entry_filter = "movie" # Use: "short" ||"movie" || "tvShort" || "tvSeries" || "tvMovie" || "tvEpisode" || "tvMiniSeries" || "tvSpecial" || "video" || "videoGame"
22
23 # Mutate data for "Average_Rating_By_Genre" and "Average_Rating_Vs_Number_Of_Films_Per_Year"
24 movies_with_rating = transmute(
25     ImdbTitleBasicsAndRatings,
26     year = startYear,
27     genre = genres,
28     titleType = titleType,
29     rating = averageRating,
30     votes = numVotes
31 ) %>% filter(titleType == data_entry_filter) %>%
32     separate_rows(genre) %>% filter(genre %in% genre_filter) %>%
33     drop_na() %>%
34     count()
35 movies_with_stats = transmute(
36     movies_with_rating,
37     year = year,
38     genre = genre,
39     # Calculate ratios to get relativ propotions of data (normalized to 0.x of 1)
40     freq_ratio = ave(freq, year, genre, FUN=sum)/ave(freq, year, FUN=sum),
41     vote_ratio = ave(votes, year, genre, FUN=sum)/ave(votes, year, FUN=sum),
42     rating_avg = ave(rating, year, genre, FUN=mean)
43 ) %>%
44     count()
45 # Mutate data for "Distribution_Of_Genres"
46 movies_without_rating = transmute(
47     ImdbTitleBasics,
48     year = startYear,
49     genre = genres,
50     titleType = titleType,
51 ) %>% filter(titleType == data_entry_filter) %>%
52     separate_rows(genre) %>% filter(genre %in% genre_filter) %>%
53     drop_na() %>%
54     count()
```

## 2. Distribution of Movie Genres

In the first plot that we generate and animate, the animation runs along the X-axis and shows the progression of the number of movies in the respective genres over the years. From 2005 onwards, almost no data seems to have been collected by the originator of the data set. To make sure that our assumption is correct, we have made a request to the data originator, so far without a reply.



From line 72 to 78, we generate the line-plot using ggplot with the years on the X-axis (x=year) and the frequency on the Y-axis (y=freq). Further on the following lines, we add titles to each of the axes and the plot.

```
72   ## Distribution_Of_Genres -- Plot ##
73   Distribution_Of_Genres = ggplot(movies_without_rating, aes(x=year, y=freq, group=genre, color=genre)) +
74     geom_line() +
75     geom_point() +
76     ggtitle("Distribution of movie genres over the past 100 years") +
77     ylab("Number of films in category") +
78     xlab("Year")
```

Then we go on with the transitions. First, we define the animation progression in relation to the year. Then we add the required parameters such as duration, frames per second, width and height and save the animation as a GIF in the project directory.
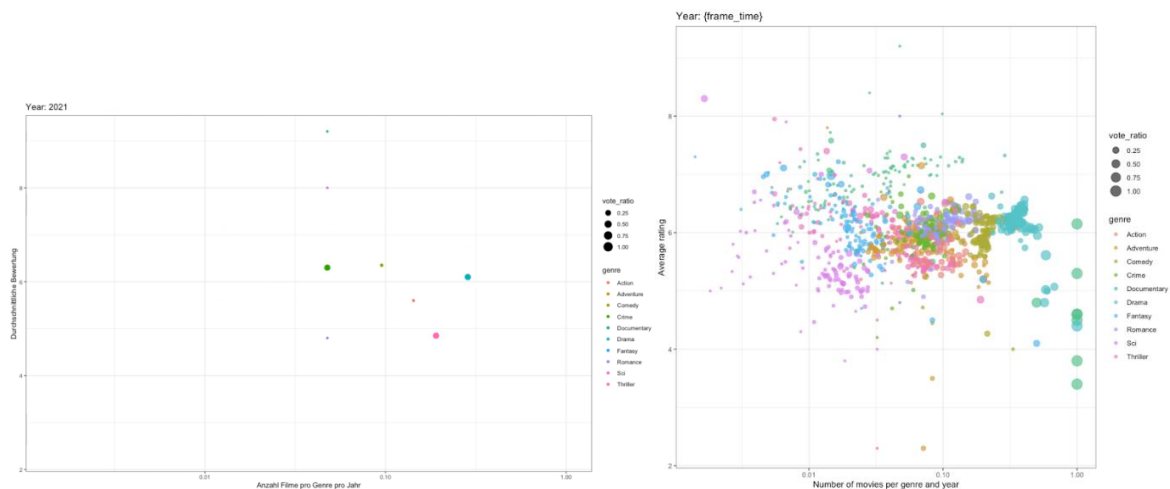
```
91   Animation_Distribution_Of_Genres = Distribution_Of_Genres +
92     transition_reveal(year) +
93     enter_fade() +
94     exit_shrink()

101  animate(Animation_Distribution_Of_Genres, duration = 10, fps = 20, width = 1000, height = 600, renderer = gifski_renderer())
102  anim_save("Distribution_Of_Genres.gif")
```

# 3. Average movie rating in Genres compared to the average number of movies produced within one year

The second plot, again from the ggplot library, is of the type geom_point, which results in the following plot. The animation represents a time progression and the points, which represent the genres each color one, vary in size depending on the proportion of votes.

Since the animation only shows the data per year (left picture), we attached the picture on the right, which shows the data of all the years, for a better understanding.



This scatterplot is generated with ggplot with a somewhat more complex axis description. We define here that the frequency of the films should be represented and that the average rating should vary in size. In addition, the genres should be displayed in different colors.

```
63   ## Average_Rating_Vs_Number_Of_Films_Per_Year -- Plot ##
64   Average_Rating_Vs_Number_Of_Films_Per_Year = ggplot(movies_with_stats, aes(freq_ratio, rating_avg, size = vote_ratio, color = genre)) +
65     #geom_point() +
66     geom_point(aes(size = vote_ratio), alpha = 0.7) +
67     scale_x_log10() +
68     theme_bw() +
69     # gganimate specific bits
70     labs(title = 'Year: {frame_time}', x = 'Number of movies per genre and year', y = 'Average rating')
```

In the following snippet we define the animation where the transitions run according to the year and we enter the data with *enter_fade()* and end by exiting the data with *exit_shrink()*.
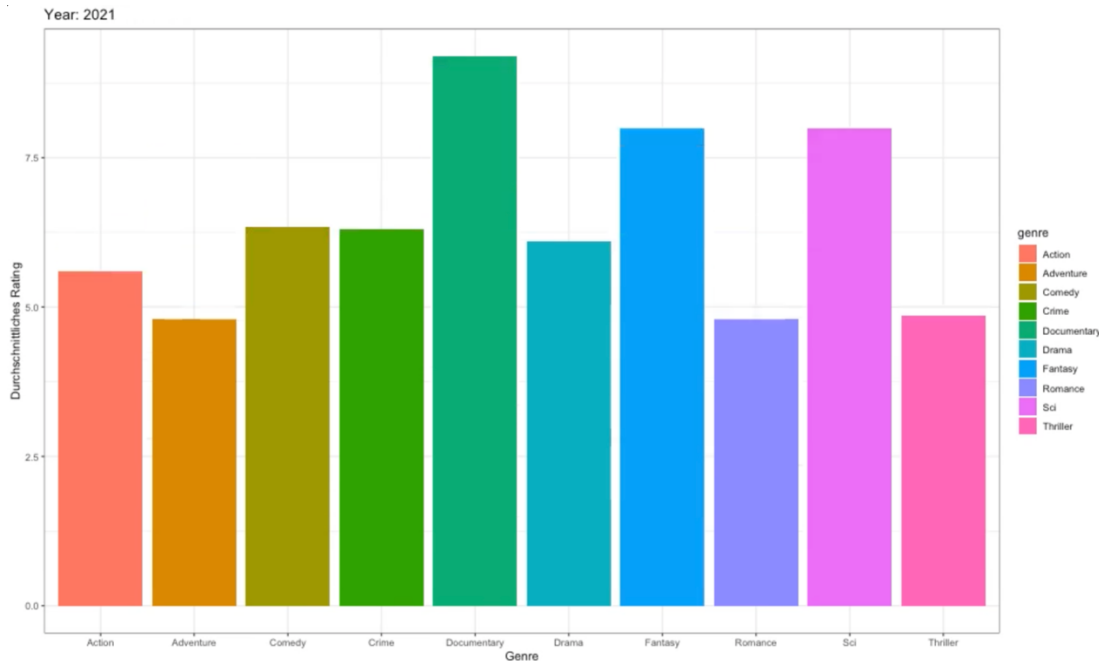
```
86   Animation_Average_Rating_Vs_Number_Of_Films_Per_Year = Average_Rating_Vs_Number_Of_Films_Per_Year +
87     transition_time(year) +
88     ease_aes('linear') +
89     enter_fade() +
90     exit_shrink()
```

Here we generate and render the gif with gifski, give a name to the output file which will be saved in the project root folder.

```
99    animate(Animation_Average_Rating_Vs_Number_Of_Films_Per_Year, duration = 30, fps = 20, width = 1000, height = 600, renderer = gifski_renderer())
100   anim_save("Average_Rating_Vs_Number_Of_Films_Per_Year.gif")
```

## 4. Average movie rating by Genres from 1897 to 2021

The third plot is a bar chart. Here the animation runs within the bars and represents a progression over more than 100 years. Shown are the average ratings (Y-axis) for the respective genres (X-axis) from 1897 to 2021.



Again, we generate the bar-plot using ggplot with the genre on the X-axis (x=genre) and the average rating on the Y-axis (y=rating_avg). We use the classic ggplot2 theme *theme_bw* and on line 61 we name the titles of each axis as well as the plot title, which will be the variable frame time in the animation, so the year.

```
56    ## Average_Rating_By_Genre -- Plot ##
57    Average_Rating_By_Genre = ggplot(movies_with_stats, aes(x=genre, y=rating_avg, fill=genre)) +
58      geom_bar(stat='identity') +
59      theme_bw() +
60      # gganimate specific bits
61      labs(title = 'Year: {frame_time}', x = 'Genre', y = 'Average rating')
```

From line 81 to 85, we define and create the animation. The values within the animation change according to the year. And again, we enter and exit the data with the known commands.

```
81    Animation_Average_Rating_By_Genre = Average_Rating_By_Genre +
82      transition_time(year) +
83      ease_aes('sine-in-out') +
84      enter_fade() +
85      exit_shrink()
```

And just like the previous plots, for the third plot, we generate an animated gif with the same animation attributes rendered with gifski and we save it with the mentioned file name in line 98 into the source folder.

```
97    animate(Animation_Average_Rating_By_Genre, duration = 20, fps = 20, width = 1000, height = 600, renderer = gifski_renderer())
98    anim_save("Average_Rating_By_Genre.gif")
```

# 5. Sources

## 5.1.    Dataset

«IMDb Dataset – From 1888 to 2023» (7[th] of April, 2021)

    Opened 12[th] of April 2021,

    From https://www.kaggle.com/komalkhetlani/imdb-dataset

## 5.2.    Chart Animations and R Code

«the R Graph Gallery - Animation» (2018)

    Opened 12[th] of April 2021,

    From https://www.r-graph-gallery.com/animation.html

«gganimate: How to Create Plots with Beautiful Animation in R» (2018)

    Opened 12[th] of April 2021,

    From https://www.datanovia.com/en/blog/gganimate-how-to-create-plots-with-beautiful-animation-in-r/

# 6. Declaration

We confirm that we have written this work independently. All passages in the text that do not originate from us are marked as quotations and provided with an exact reference to their origin.

The sources used (also applies to illustrations, graphics, etc.) are listed in the bibliography or list of sources.

Birnbaumer Jamie-Lee          Djerrah Yacine          Zimmermann Yannik