

# Statistical Consulting Final Write Up

Xunqing Zheng & Duja Michael

## Extracting and Analyzing the New York Times Movie Reviews Data

### Introduction

As the movie landscape continues to evolve, The New York Times movie reviews continue to offer insightful commentary, providing audiences with a nuanced perspective on the ever-changing and dynamic world of cinema. In order to get a deeper understanding of how the NYTimes is selecting what movies to review, what type of movie each of the staff critics review, and if there are patterns selecting which movies to be reviewed by whom, we decided to analyze all the movie reviews published in NYTimes from January 2018 until November 2023. In particular, we are trying to understand:

- Is there a pattern or a time-trend in the number of movies that the NYTimes reviews during this period?
- Is any particular genre more likely to be reviewed by NYTimes?
- Which critics are the most prolific?
- Do particular reviewers/authors specialize in particular genres?
- What's the distribution of IMDb ratings of the movies reviewed each author?
- Are "Critic's Picks" more likely to also receive higher ratings by the public?

### Initial Exploration

Our journey to land on this topic was a windy one. Initially, we wanted to work with a social media API where we would grab all the comments on a specific movie post and run sentiment analysis to get an idea about how the movie is received by the public. We also had a detour in our thinking in mid-October and wanted to shift our project to check the state of democracy in the USA by downloading the comments on all posts by the POTUS and Bernie Sanders accounts to see the prevalence of the calls for a ceasefire in Gaza knowing that those comments went on deaf ears. All those ideas proved to be unfit for the scope of this project due to the difficulty in extracting comments on a certain post in the APIs of the social media platforms we wanted to use.

The final pivot we did was to go back to our initial idea of movie reviews and use the NYTimes' API.

### Difficulties with the NYTimes API

The NYTimes API presents a steep learning curve on how it operates. A good understanding of the documentation is required to navigate and extract the desired information efficiently. Throughout the process, we have encountered many difficulties, such as function deprecation, rate limits, data consistency, etc.

#### API Limitations:

The NYTimes API has a daily rate limit of 500 hits per API key, each hit only gets 10 articles, and it also limits users from hitting the API too frequently. An error message will appear once we reached 200

consecutive request to access the API, sometimes even earlier. These limitations impacted the real time data retrieval and require careful consideration of the API usage patterns. On top of the rate limits, NYTimes also deprecated its movie reviews API, which made our data collection process even more time consuming. We used its article search API as a workaround. The issue with the article search API is that in order to get all the movie reviews, we can't limit the articles to only movie reviews. A lot of other articles were included in our pull and only 28% of the article were actually movie reviews.

### Code Errors:

Working with API is never easy, in addition to the API limitation, we also experienced some errors in our codes, which further complicates our data collection process. Although NYTimes deprecated its Movie Review API, it offered a sample url structure to access the movie review from the Article Search API.

- [https://api.nytimes.com/svc/search/v2/articlesearch.json?fq=section\\_name%3A%22Movies%22&type\\_of\\_material%3A%22Reviews%22](https://api.nytimes.com/svc/search/v2/articlesearch.json?fq=section_name%3A%22Movies%22&type_of_material%3A%22Reviews%22)  
key{your-key}

### Quality, consistency, and insufficient information:

Another issue we faced with the NYTimes API was the lack of information and consistency in the data provided. For instance, data on the movie name and movie type are provided in a list of lists. Not all reviews had these two pieces of information, and in the cases where the information is provided, it is not arranged in a way that lends itself to systematic extraction. For movie names, we had to supplement the list data with information from the URL in cases where the list was missing that information. For the type of movie, the only categorizations that were consistently provided were whether the movie is a film, a documentary, or an animated film. No consistent information is provided for the genre of the movie.

From the NYTimes API data, we ended up extracting and using movie name, movie type, whether or not it was a critic's pick, the author of the review, and the date of publishing the review.

## Supplementing the NYTimes data with IMDb information

Besides its paid API, IMDb has free datasets that are available for use for non-commercial purposes. We utilize this data to supplement the NYTimes information. In particular, this data adds information on the genre of the movie, the IMDb average rating, and the number of ratings.

## Data Collection

**Xunqing** → talk about how we got the info from the NYT api and how we saved it for later analysis;

**unqing** → Also talk about the chunk of code you created to get the IMDb data directly from the web and about the two datasets (title basics and title reviews)

```
imdbTSVfiles <- function(fileName){  
  url <- paste0("https://datasets.imdbws.com/",fileName,".tsv.gz")  
  tmp <- tempfile()  
  download.file(url, tmp)  
  
  assign(fileName,  
    readr::read_tsv(  
      file = gzfile(tmp),
```

```

      col_names = TRUE,
      quote = "",
      na = "\\N"),
      envir = .GlobalEnv)
}

```

```
imdbTSVfiles("title.basics")
```

```

## Rows: 10415898 Columns: 9
## -- Column specification -----
## Delimiter: "\t"
## chr (5): tconst, titleType, primaryTitle, originalTitle, genres
## dbl (4): isAdult, startYear, endYear, runtimeMinutes
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

```

```
imdbTSVfiles("title.ratings")
```

```

## Rows: 1381661 Columns: 3
## -- Column specification -----
## Delimiter: "\t"
## chr (1): tconst
## dbl (2): averageRating, numVotes
##
## i Use 'spec()' to retrieve the full column specification for this data.
## i Specify the column types or set 'show_col_types = FALSE' to quiet this message.

```

## NYTimes data manipulation steps

1. All the files that were extracted from the API and saved as .R files were loaded and appended using the following code:

```

library(tidyverse)
library(lubridate)
library(stringr)
library(zoo)
library(readr)
library(tidyverse)
library(splitstackshape)
library(stringi)
library(tm)
library(ggplot2)

```

```

#####
# loading API movie data
#####

# appending the data
Movie2018H1 <- load("Movie2018H1.Rdata")
Movie2018H1 <- MovieFile

```

```

Movie2018H2 <- load("Movie2018H2.Rdata")
Movie2018H2 <- MovieFile
Movie2019H1 <- load("Movie2019H1.Rdata")
Movie2019H1 <- MovieFile
Movie2019H2 <- load("Movie2019H2.Rdata")
Movie2019H2 <- MovieFile
Movie2020H1 <- load("Movie2020H1.Rdata")
Movie2020H1 <- MovieFile
Movie2020H2 <- load("Movie2020H2.Rdata")
Movie2020H2 <- MovieFile
load("Movie2021H1.Rdata")
load("Movie2021H2.Rdata")
load("Movie2022H1.Rdata")
load("Movie2022H2.Rdata")
load("Movie2023H1.Rdata")
load("Movie2023H2.Rdata")

```

```

merged_all <- rbind(Movie2018H1, Movie2018H2,
                    Movie2019H1, Movie2019H2,
                    Movie2020H1, Movie2020H2,
                    Movie2021H1, Movie2021H2,
                    Movie2022H1, Movie2022H2,
                    Movie2023H1, Movie2023H2)

```

2. Given that the API provides movie and book reviews, we filter out all the observations that are not movie reviews using information from the URL column:

```

# keeping only movies using the /movies/ pattern in response.docs.web_url
merged <- merged_all %>%
  filter(str_detect(response.docs.web_url, "/movies"))
# we have 4268 movies, now keep only unique ones

```

This gives a dataset with 4268 rows (movie reviews).

3. Only keep the unique reviews since some of the reviews might have been pulled more than once due to an overlap of the time period in the API call. This leaves us with 4142 unique movie reviews.

```
merged_unique <- merged %>% distinct(response.docs.snippet, .keep_all = TRUE)
```

4. Create columns with variables of interest including review author, critic's pick, review date, movie type, and movie name:

```

# Add author name (removing "By")
stopwords = c("By")
merged_unique$author <-gsub("By","",as.character(merged_unique$response.docs.byline.original))

# clean critic's pick column
merged_unique$criticpick <- ifelse(is.na(merged_unique$response.docs.headline.kicker),0,1)

# clean date column

```

```
merged_unique <- merged_unique %>%
  mutate (month = month(response.docs.pub_date),
          year = year(response.docs.pub_date))
merged_unique$monthyear <- as.yearmon(paste(merged_unique$year, merged_unique$month), "%Y %m")

# loop to extract movie type and movie name
merged_unique <- merged_unique %>% mutate(type=0,
                                          movie_name=0)

for (i in 1:nrow(merged_unique)) {
  typetable=as.data.frame(merged_unique$response.docs.keywords[i])
  merged_unique$type[i]=ifelse(any(typetable$value=="Documentary Films and Programs"),"Documentary Films",
                              ifelse(any(typetable$value=="Movies"),"Movies",
                                      ifelse(any(typetable$value=="Animated Films"),"Animated Films",NA)))
  merged_unique$movie_name1[i] <- typetable[typetable$name=="creative_works", "value"][1]
}

table(merged_unique$type, useNA = "always") #12 movies have no type information
```

```
##
##           Animated Films Documentary Films and Programs
##           13                                1095
##           Movies                                <NA>
##           3022                                12
```

```
#Extract movie name from URL
merged_unique$movie_name2 <- str_match(merged_unique$response.docs.web_url, "/movies/\\s*(.*?)\\s*-review")

#Use name from URL if the information is not there in the response.docs.keywords list
merged_unique$movie_name <- ifelse(is.na(merged_unique$movie_name1),merged_unique$movie_name2,merged_unique$movie_name1)

check3 <- merged_unique %>% filter(is.na(movie_name)) #only 3 have unidentifiable name from either source
```

5. Clean up the movie names sufficiently so that they can be effectively merged with the IMDb data

```
# remove (Movie) ; replace dashes with spaces ; remove white space ; lower case all for movie names
merged_unique$movie_name <-str_replace(merged_unique$movie_name, " \\s*\\\[^\[\\]\]+\\)", "")
merged_unique$movie_name <-str_replace_all(merged_unique$movie_name,'-', ' ')
merged_unique$movie_name <- trimws(merged_unique$movie_name)
merged_unique$movie_name <- tolower(merged_unique$movie_name)
# replace & with "and", remove ":"
merged_unique$movie_name <-str_replace_all(merged_unique$movie_name,':','')
merged_unique$movie_name <-str_replace_all(merged_unique$movie_name,'&','and')
```

6. Keep relevant columns only and get the number of reviews per month and year

```
api_unique <- merged_unique %>% dplyr::select (response.docs.pub_date, month, year, monthyear,
                                             response.docs.news_desk, type, author, criticpick,
                                             movie_name, movie_name1, movie_name2) %>%

  group_by(monthyear) %>%
  mutate(n_month=n()) %>%
  ungroup() %>% group_by(year) %>%
```

```
mutate(n_year=n()) %>%
ungroup() %>%
group_by(monthyear) %>%
arrange(monthyear)
```

## IMDb data preparation

1. Keep movies from 2017 onwards only since the NYTimes movie reviews start in 2017 and it is very unlikely for the NYTimes to review movies from over a year before and merge this information with that from the ratings dataset by movie ID (tconst).

```
titles <- title.basics %>%
  filter(startYear>2016,
         titleType=="movie"|titleType=="tvMovie")
ratings <- title.ratings
joined <- left_join(titles,ratings, by="tconst") %>%
  rename(movie_name = primaryTitle,
         release_year=startYear) %>%
  mutate(release_year=as.numeric(release_year))
```

2. Clean up movie names in preparation for the merge with NYTimes data

```
# Clean up movie names: replace dashes with spaces ; remove white space ; lower case all for movie name.
joined$movie_name <-str_replace_all(joined$movie_name,'-', ' ')
joined$movie_name <- trimws(joined$movie_name)
joined$movie_name <- tolower(joined$movie_name)
joined$movie_name <-str_replace_all(joined$movie_name,':','')
joined$movie_name <-str_replace_all(joined$movie_name,'&','and')
```

## Merging the two data sources

The code chunk below shows the merge and the filtering out of non-unique matches from the two data sources

```
withratings <- left_join(api_unique, joined, by="movie_name") #5350
```

```
## Warning in left_join(api_unique, joined, by = "movie_name"): Detected an unexpected many-to-many relationship.
## i Row 3 of 'x' matches multiple rows in 'y'.
## i Row 47319 of 'y' matches multiple rows in 'x'.
## i If a many-to-many relationship is expected, set 'relationship =
## "many-to-many"' to silence this warning.
```

```
# keep only the joins where release date is BEFORE review date
withratings <- withratings %>% filter(release_year <= year) %>%
  group_by(movie_name) %>% mutate(n_matches=n(),
                                unclear=ifelse(n_matches>1,1,0)) #4456 row left
table(withratings$n_matches, useNA = "always") # 3232 uniquely matched out of 4142
```

```
##
##      1      2      3      4      5      6      7      8      9     10     12 <NA>
## 3235  516  264  156  115   48   49   16   27   10   12    0
```

3232/4142

```
## [1] 0.7802994
```

```
#keep only matched movies  
matched <- withratings %>% filter(n_matches==1)
```

3232 out of 4142 movies were uniquely matched with the IMDb data. A random sub-sample of those were manually checked, and they were all true matches. Our match rate is roughly 80%, which is not bad.

## Transforming genre data to a usable format

```
# get the genres:  
gens <- unique(word(matched$genres, sep=","))  
length(gens) #19 genres represented
```

```
## [1] 19
```

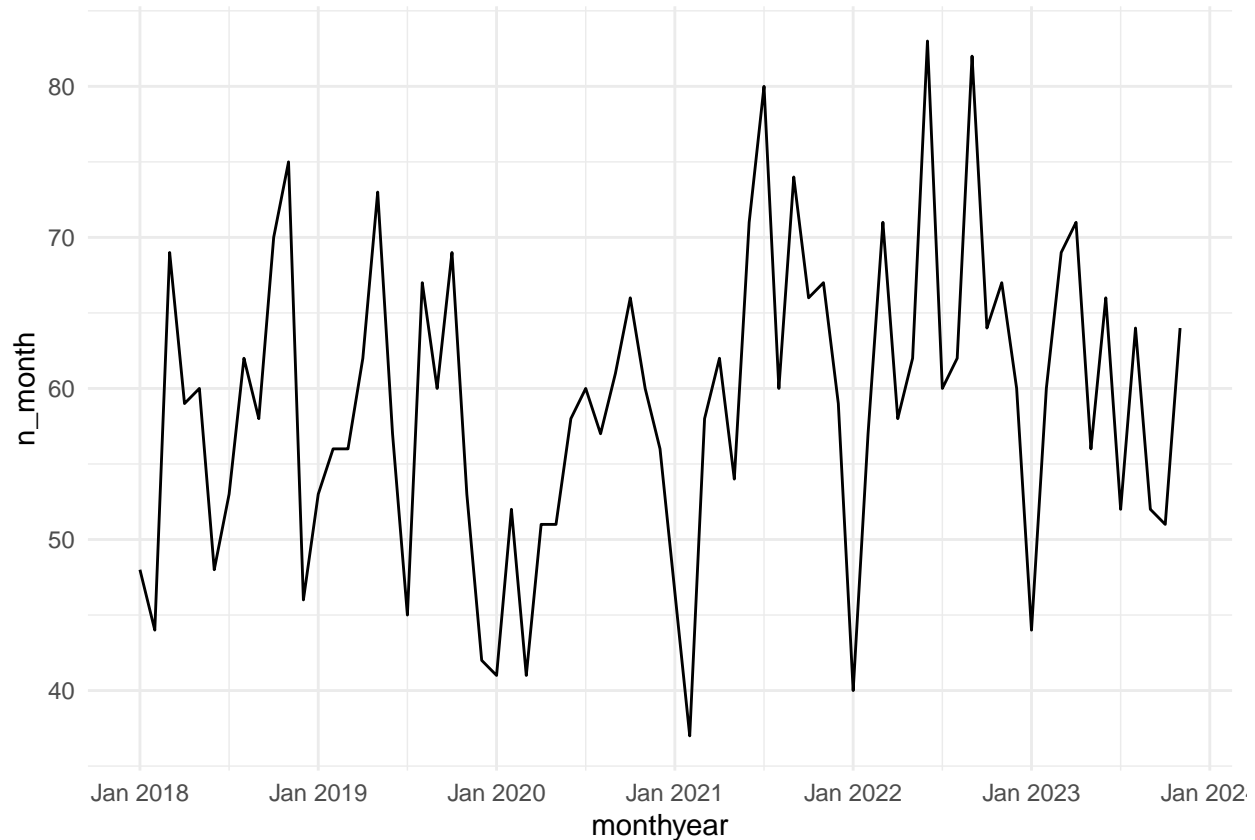
```
# make them into indicator columns:  
matched$Documentary <- grepl("Documentary", matched$genres)  
matched$Drama <- grepl("Drama", matched$genres)  
matched$Animation <- grepl("Animation", matched$genres)  
matched$Comedy <- grepl("Comedy", matched$genres)  
matched$Crime <- grepl("Crime", matched$genres)  
matched$Action <- grepl("Action", matched$genres)  
matched$Adventure <- grepl("Adventure", matched$genres)  
matched$Biography <- grepl("Biography", matched$genres)  
matched$Horror <- grepl("Horror", matched$genres)  
matched$Mystery <- grepl("Mystery", matched$genres)  
matched$Animation <- grepl("Animation", matched$genres)  
matched$Thriller <- grepl("Thriller", matched$genres)  
matched$SciFi <- grepl("Sci-Fi", matched$genres)  
matched$Fantasy <- grepl("Fantasy", matched$genres)  
matched$Family <- grepl("Family", matched$genres)  
matched$Musical <- grepl("Musical", matched$genres)  
matched$History <- grepl("History", matched$genres)  
matched$Music <- grepl("Music", matched$genres)  
matched$Romance <- grepl("Romance", matched$genres)
```

## Analysis and Visualizations

```
#review trends  
api_unique %>%  
  arrange(monthyear) %>%  
  group_by(monthyear, year) %>%  
  summarise(n_month=n()) %>%  
  ggplot() +  
  aes(x = monthyear, y = n_month) +
```

```
geom_line()+
scale_color_gradient() +
theme_minimal()
```

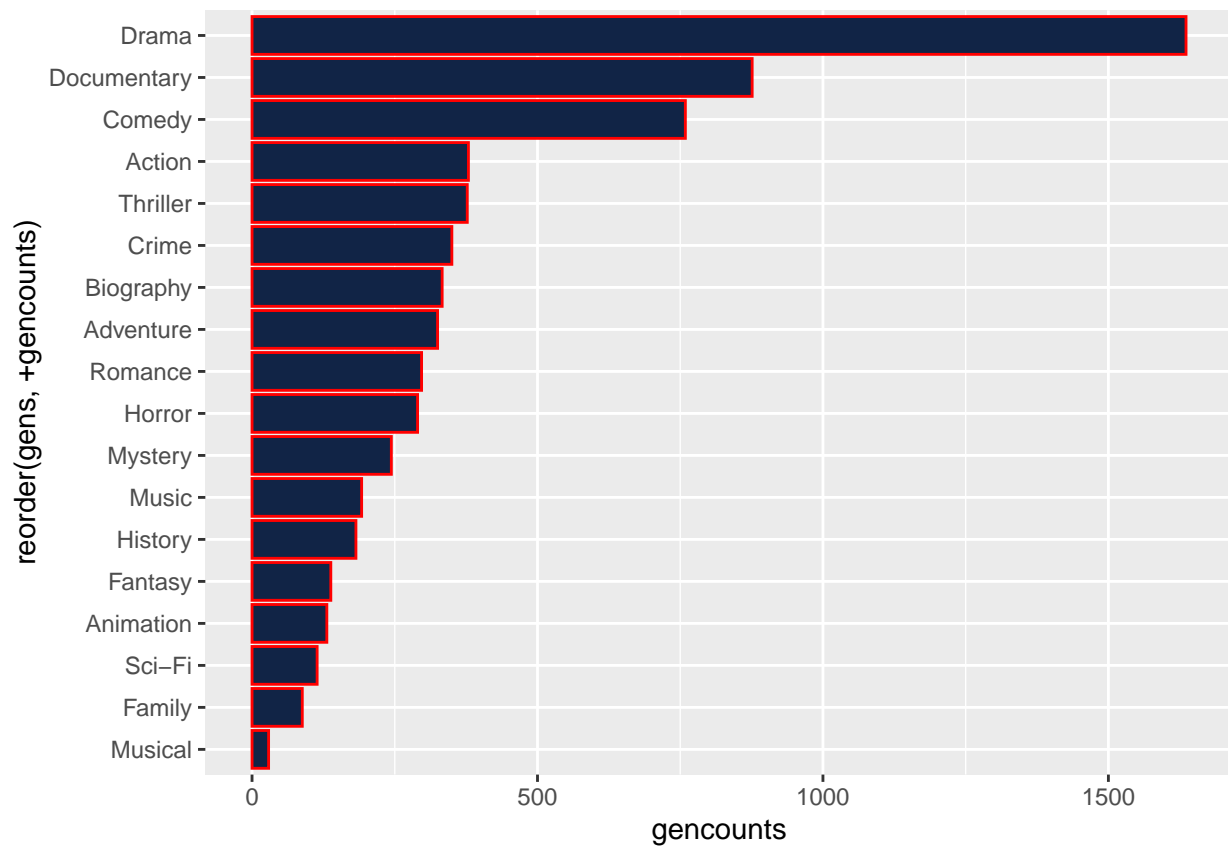
## 'summarise()' has grouped output by 'monthyear'. You can override using the  
## '.groups' argument.



```
#distribution of reviewed movies across genres
gens <- gens[-14]
gencounts <- c(sum(matched$Documentary==1), sum(matched$Drama==1),
               sum(matched$Comedy==1), sum(matched$Crime==1),
               sum(matched$Action==1), sum(matched$Adventure==1),
               sum(matched$Biography==1), sum(matched$Horror==1),
               sum(matched$Mystery==1), sum(matched$Animation==1),
               sum(matched$Thriller==1), sum(matched$SciFi==1),
               sum(matched$Fantasy==1), sum(matched$Family==1),
               sum(matched$Musical==1), sum(matched$History==1),
               sum(matched$Music==1), sum(matched$Romance==1))
genre_table=as.data.frame(cbind(gens,gencounts))
genre_table$gencounts <- as.numeric(genre_table$gencounts)
genre_table <- genre_table[order(genre_table$gencounts, decreasing = TRUE), ]

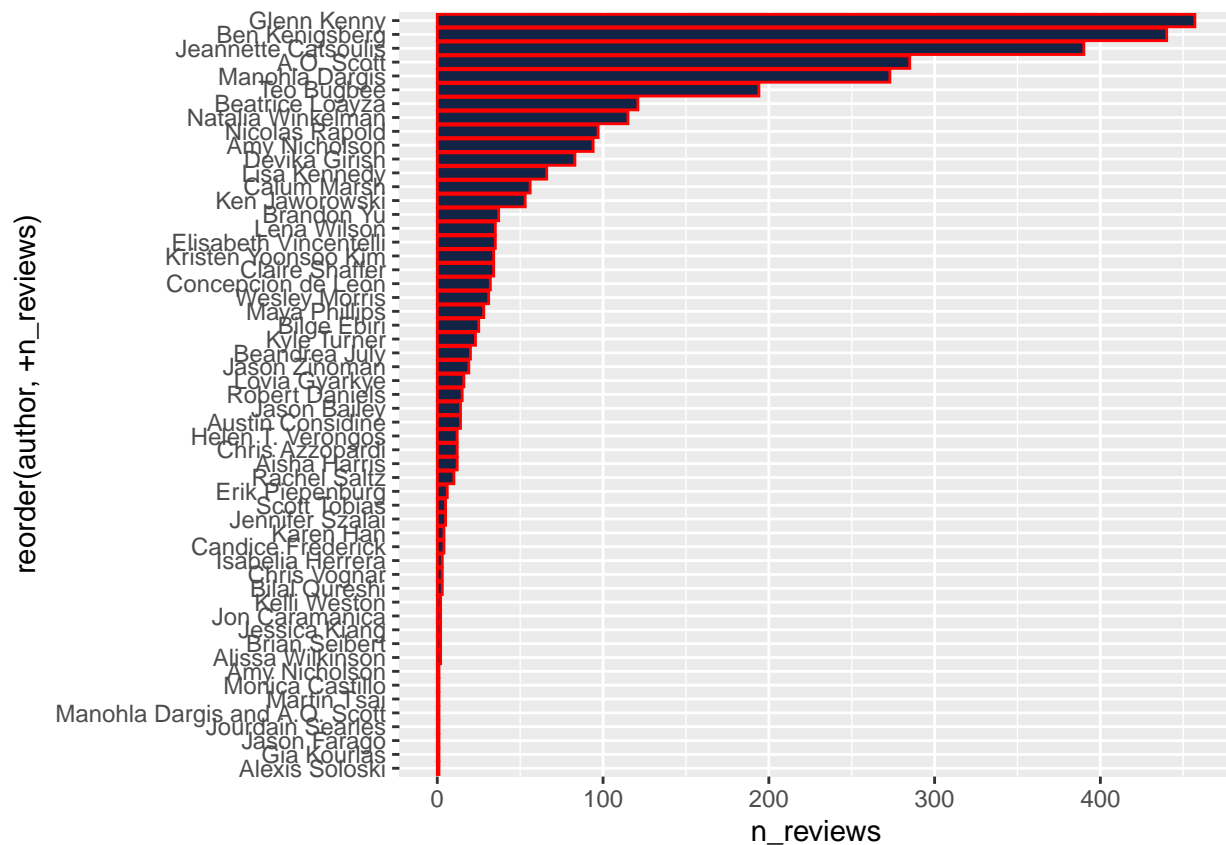
p <- ggplot(genre_table, aes(x = reorder(gens, +gencounts), y = gencounts)) +
  geom_bar(stat="identity", color='red',fill="#112446") +
  coord_flip()
p
```





```
#number of reviews per author
auth_rate <- matched %>% group_by(author) %>%
  summarise(mean_imdb_rating=mean(averageRating, na.rm=TRUE),
            n_reviews=n())
auth_rate <- na.omit(auth_rate)

p <- ggplot(auth_rate, aes(x = reorder(author, +n_reviews), y = n_reviews)) +
  geom_bar(stat="identity", color='red', fill="#112446") +
  coord_flip()
p
```



*#genre per author*

```
genres_author <- matched %>% group_by(author) %>% summarise(n=n(),
  perc_pick=mean(criticpick),
  perc_rating <- mean(averageRating, na.rm=TRUE),
  perc_Documentary=mean(Documentary),
  perc_Drama=mean(Drama),
  perc_Animation=mean(Animation),
  perc_Comedy=mean(Comedy),
  perc_Crime=mean(Crime),
  perc_Action=mean(Action),
  perc_Adventure=mean(Adventure),
  perc_Biography=mean(Biography),
  perc_Horror=mean(Horror),
  perc_Mystery=mean(Mystery),
  perc_Thriller=mean(Thriller),
  perc_SciFi=mean(SciFi),
  perc_Fantasy=mean(Fantasy),
  perc_Family=mean(Family),
  perc_Musical=mean(Musical),
  perc_History=mean(History),
  perc_Music=mean(Music),
  perc_Romance=mean(Romance)) %>%

  filter(n>9)
na.omit(genres_author)
```

## # A tibble: 34 x 22

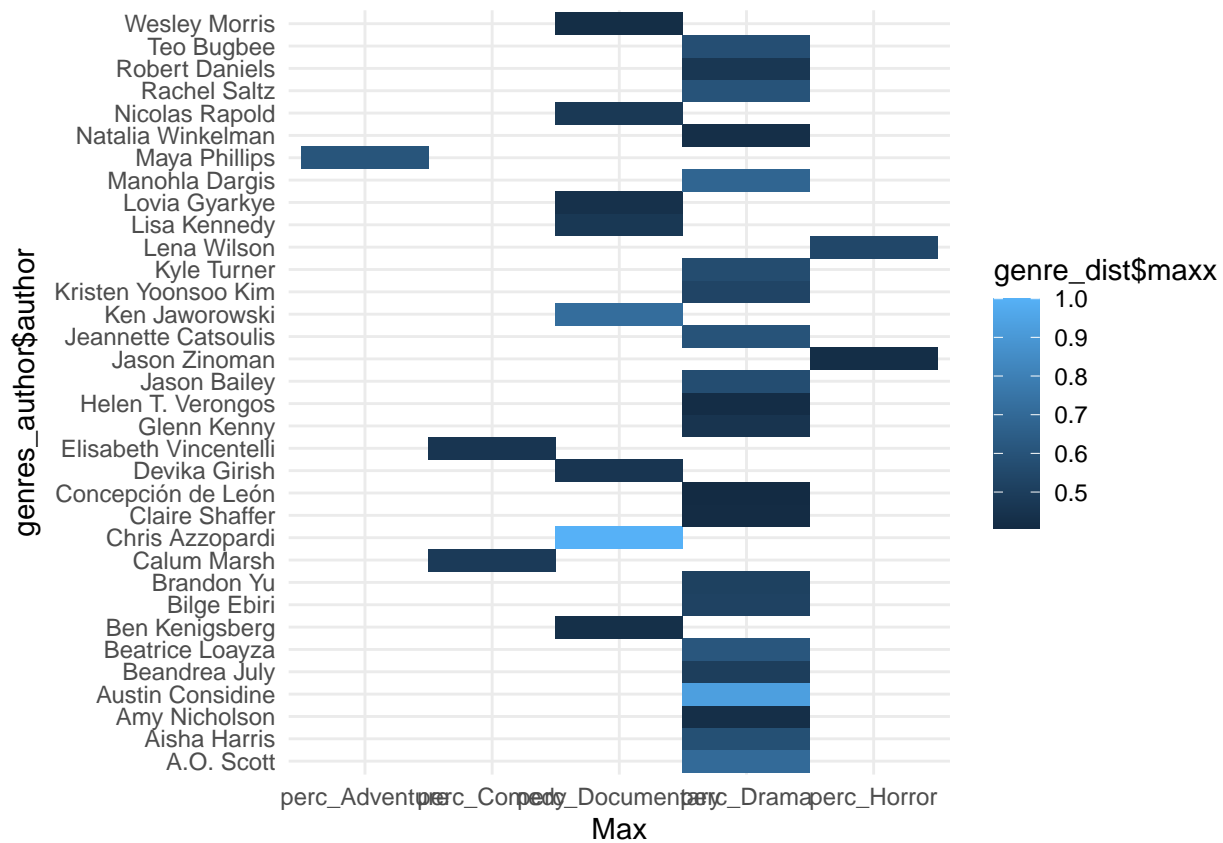
```
##   author          n perc_pick perc_rating <- mean(~1 perc_Documentary perc_Drama
##   <chr>          <int>      <dbl>          <dbl>          <dbl>          <dbl>
## 1 " A.O. Sc~    285      0.291          6.78          0.119          0.702
## 2 " Aisha H~    12       0.167          6.58          0.0833         0.583
## 3 " Amy Nic~    94       0.213          6.23          0.0957         0.426
## 4 " Austin ~    14       0.286          6.34          0           0.929
## 5 " Beandre~    20       0.25           6.42          0.35           0.5
## 6 " Beatric~   121       0.149          6.50          0.190          0.612
## 7 " Ben Ken~   440       0.114          6.55          0.430          0.395
## 8 " Bilge E~    25       0.28           6.31          0.24           0.52
## 9 " Brandon~   37       0.162          6.21          0.162          0.514
## 10 " Calum M~   56       0.0893         6.1           0.214          0.268
## # i 24 more rows
## # i abbreviated name: 1: 'perc_rating <- mean(averageRating, na.rm = TRUE)'
## # i 16 more variables: perc_Animation <dbl>, perc_Comedy <dbl>,
## #   perc_Crime <dbl>, perc_Action <dbl>, perc_Adventure <dbl>,
## #   perc_Biography <dbl>, perc_Horror <dbl>, perc_Mystery <dbl>,
## #   perc_Thriller <dbl>, perc_SciFi <dbl>, perc_Fantasy <dbl>,
## #   perc_Family <dbl>, perc_Musical <dbl>, perc_History <dbl>, ...
```

```
genre_dist <- genres_author %>%
  mutate(maxx = pmax(perc_Documentary, perc_Drama,perc_Comedy, perc_Crime,
    perc_Action,perc_Adventure,perc_Biography,perc_Horror,
    perc_Mystery,perc_Animation,perc_Thriller,perc_SciFi ,
    perc_Fantasy,perc_Family,perc_Musical ,perc_History,
    perc_Music,perc_Romance))

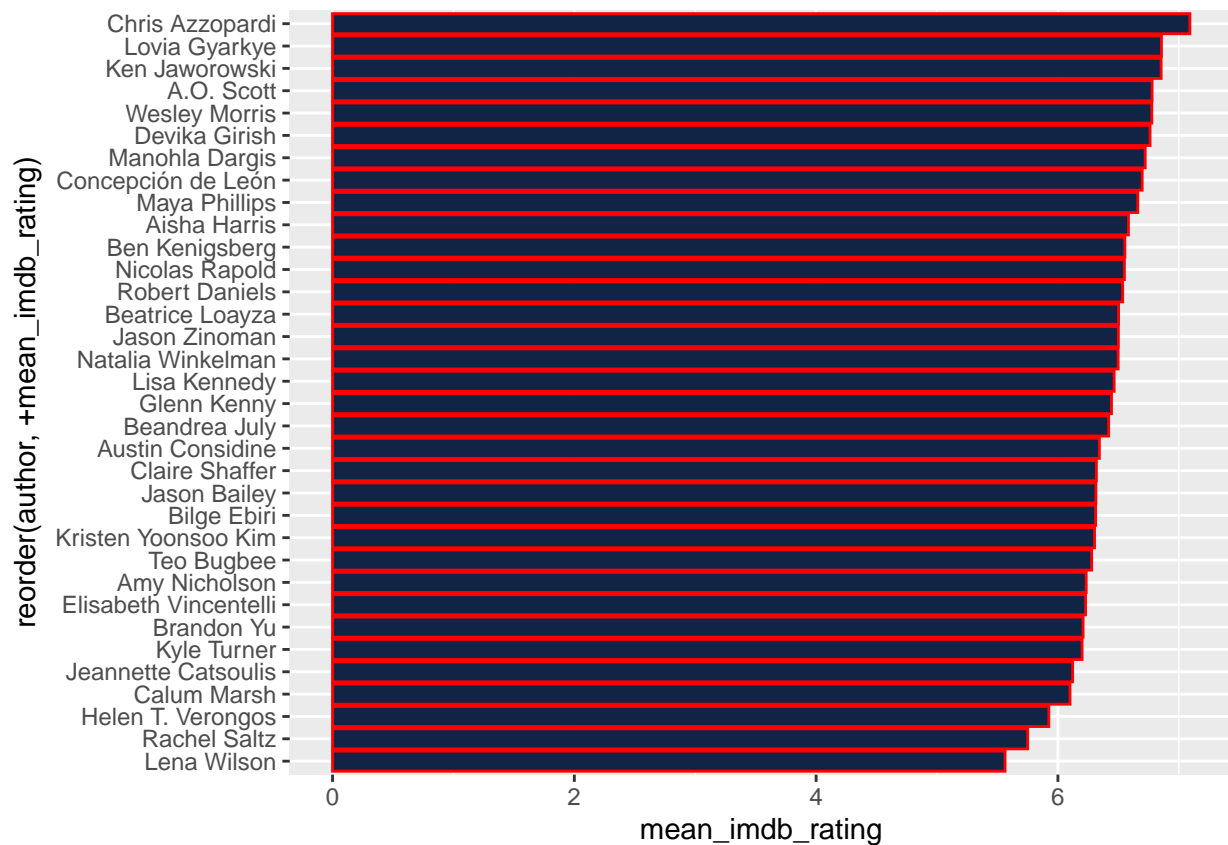
genre_max <- genres_author %>%
  dplyr::select(perc_Documentary, perc_Drama,perc_Comedy, perc_Crime,
    perc_Action,perc_Adventure,perc_Biography,perc_Horror,
    perc_Mystery,perc_Animation,perc_Thriller,perc_SciFi ,
    perc_Fantasy,perc_Family,perc_Musical ,perc_History,
    perc_Music,perc_Romance) %>%
  rowwise %>%
  mutate(Max = names(.)[which.max(c(perc_Documentary, perc_Drama,perc_Comedy, perc_Crime,
    perc_Action,perc_Adventure,perc_Biography,perc_Horror,
    perc_Mystery,perc_Animation,perc_Thriller,perc_SciFi ,
    perc_Fantasy,perc_Family,perc_Musical ,perc_History,
    perc_Music,perc_Romance))]) %>% ungroup

author_n_max <- cbind(genres_author$author,genre_max,genre_dist$maxx)

ggplot(author_n_max) +
  aes(
    x = `genres_author$author`,
    y = Max,
    fill = `genre_dist$maxx`
  ) +
  geom_tile() +
  scale_fill_gradient() +
  coord_flip() +
  theme_minimal()
```



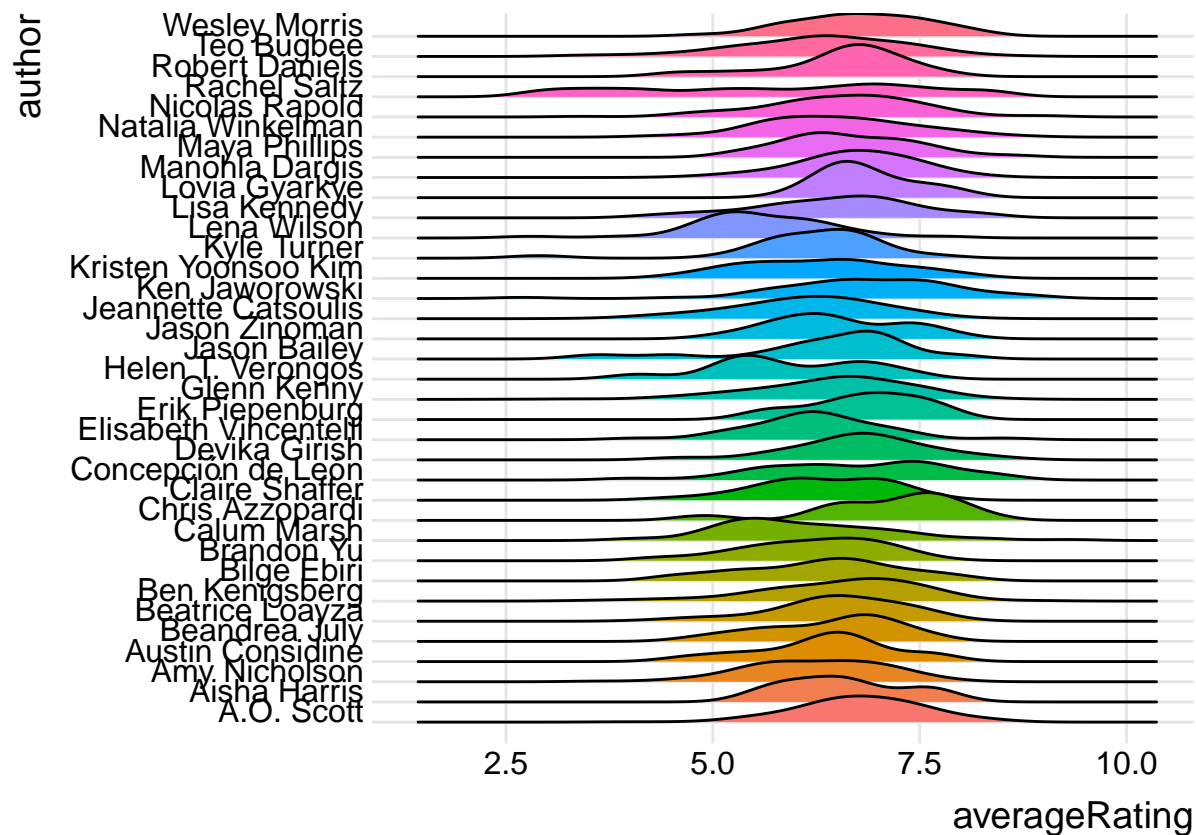
```
#average imdb rating per author (those with more than 10 reviews)
auth_rate <- matched %>% group_by(author) %>%
  summarise(mean_imdb_rating=mean(averageRating, na.rm=TRUE),
            n_reviews=n()) %>%
  filter(n_reviews>9)
auth_rate <- na.omit(auth_rate)
p <- ggplot(auth_rate, aes(x = reorder(author, +mean_imdb_rating), y = mean_imdb_rating)) +
  geom_bar(stat="identity", color='red',fill="#112446") +
  coord_flip()
p
```



```
library(ggribes)
ridge <- matched %>% group_by(author) %>% mutate(n=n()) %>% filter(n>5)
ggplot(ridge, aes(x = averageRating, y = author, fill = author)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.354
```

```
## Warning: Removed 16 rows containing non-finite values
## ('stat_density_ridges()').
```



*#critics pick ratings vs non critics pick*

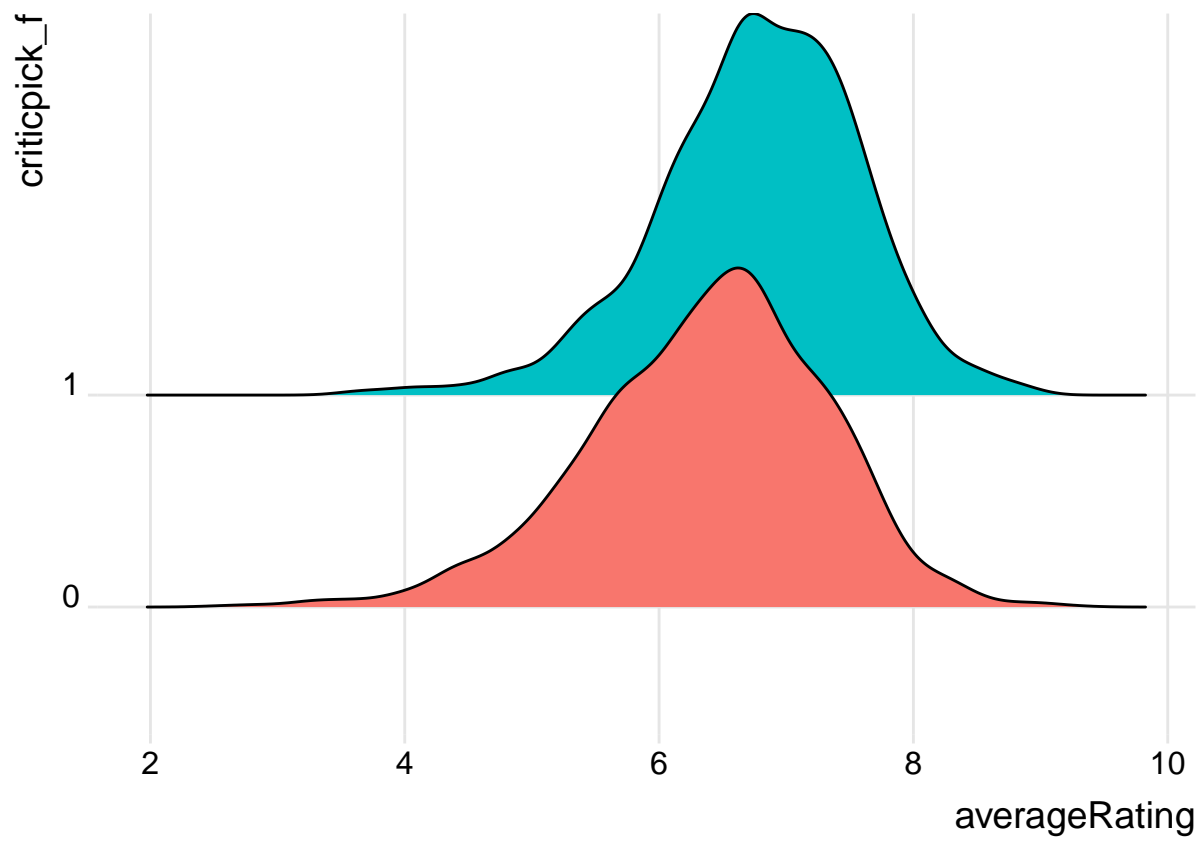
```
matched %>% group_by(criticpick) %>% summarise(avg_imdb_rating=mean(averageRating, na.rm=TRUE))
```

```
## # A tibble: 2 x 2
##   criticpick avg_imdb_rating
##   <dbl>         <dbl>
## 1       0           6.38
## 2       1           6.78
```

```
matched$criticpick_f=as.factor(matched$criticpick)
ggplot(matched, aes(x = averageRating, y = criticpick_f, fill = criticpick_f)) +
  geom_density_ridges() +
  theme_ridges() +
  theme(legend.position = "none")
```

```
## Picking joint bandwidth of 0.175
```

```
## Warning: Removed 17 rows containing non-finite values
## ('stat_density_ridges()').
```



Conclusion