

Anime Popularity Analysis

Zening Ye

11/22/2021

Abstract

In this project, I will use the data from 2014 to 2017 to analyze the following questions: first, what source(s) of Anime affected the audience more in 2014 to 2017? Second, why did the Anime come up a lot in 2014 to 2017? Last but not least, what type(s) of Anime will be the mainstream in the future? The group-level of this project will be the source of the Anime, there are 12 groups for the source. It is more reliable to fit the model. Last but not least, The report has three aspects to illustrate the entire research: Introduction, Method and Result.

Introduction

Anime has become more and more popular for the past 5 to 6 years. The original word “Anime” means a hand-drawn and computer animation originating from Japan. In Japanese and Japanese, the Anime includes all animated works, regardless of style or origin. There are plenty of Anime works in the world, such as Movie, Music, ONA, TV, etc. However, some Anime did not have more popularity than others, and they did not have more shows on TV or other.

There are plenty of types, sources and genres in Anime, therefore there will be a lot of combinations whether in TV, Movie, OVA, etc. There are some variables I would like to use to fit the model, such as start day, types, genres, sources, rating, popularity and score. These variables might help me to get some direction for the questions I mentioned above. For instance, how the genres might affect the popularity of the source(manga, original, novel), how rating and score will affect the popularity as well. Furthermore, how will these variables influence the future mainstream of the Anime?

Based on these factors, I would like to use multilevel models to illustrate how these factors affect the mainstream and popularity in single or multiple genres in the future.

Method

Data Processing and Cleaning

The data was came from Kaggle: Anime dataset:. The dataset I have chosen was cleared, however, it was not good enough for me to move to the next steps. Therefore, I first tidy up the data. I removed the variables that I will not use for the analysis in this report such as title, title_english, synopsis, etc. Second, I renamed some columns for easy access in the future, such as arid_from to star_year. Since I’m focusing on the data from 2008 to 2020, I filtered the data from the original dataset and created a new dataset for exploratory data analysis, which might help to figure out the difference from two different timelines. After cleaning the datasets I will use for the further analysis is listed below:

Column	Description
title	Name of the Anime
year_sta	Date on Aired
duration	Duration of Each Episode in min and hour
episodes	Number of episodes in Anime
genres	Theme of the Anime
popularity	Famous Level
rank	Rank of the Anime
rating	Level of the Anime, eg. PG-13, R, PG
score	Numerical Level for Anime
scored_by	Population of the rating
source	Category of the Anime
type	Type of the Anime

Exploratory Data Analysis

As I mentioned above I will use the data from 2008 to 2020, since during this period, more and more Anime came up, the popularity, rating, and sources are dynamic. The plot below indicates the frequency of the

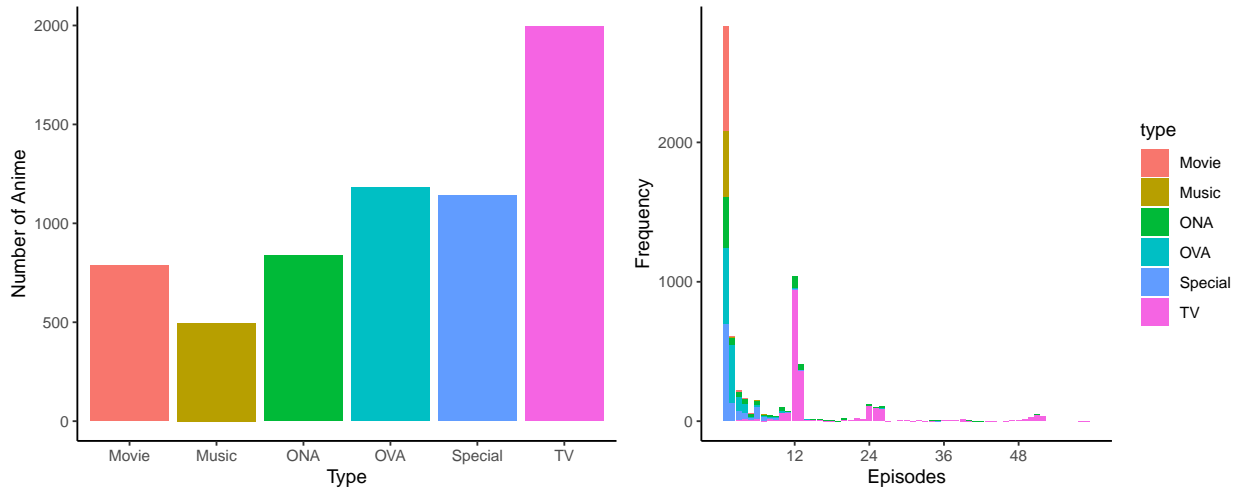


Figure 1: Frequency of Anime in 2008-2020

For Figure 1, it is obvious to see the type of Anime in TV with episodes under 24 has more popularity and market share in 2008 to 2020. In addition, the second plot in Figure 1 basically represents the entire then entire stream type of Anime in these years.

For Figure 2, in order to show the distribution of Anime in different types, I filtered the original dataset into six different subsets by different types. For instance, I used TV and Movie as examples to illustrate the frequency of source. Apparently, Original and Light Novels are the most famous topics in TV and Movie.

Model Fitting

Before I fit the multilevel model, I would like to use the subset I made above to fit a glm model to check their relationship. Since TV and Movies are the most effective sources on Anime, it might help to visualize the relationship. The variables will include score, rating, source, and rank. Since the popularity level might be too large for fitting the model, I use $\log(\text{popularity})$ to fit the linear model.

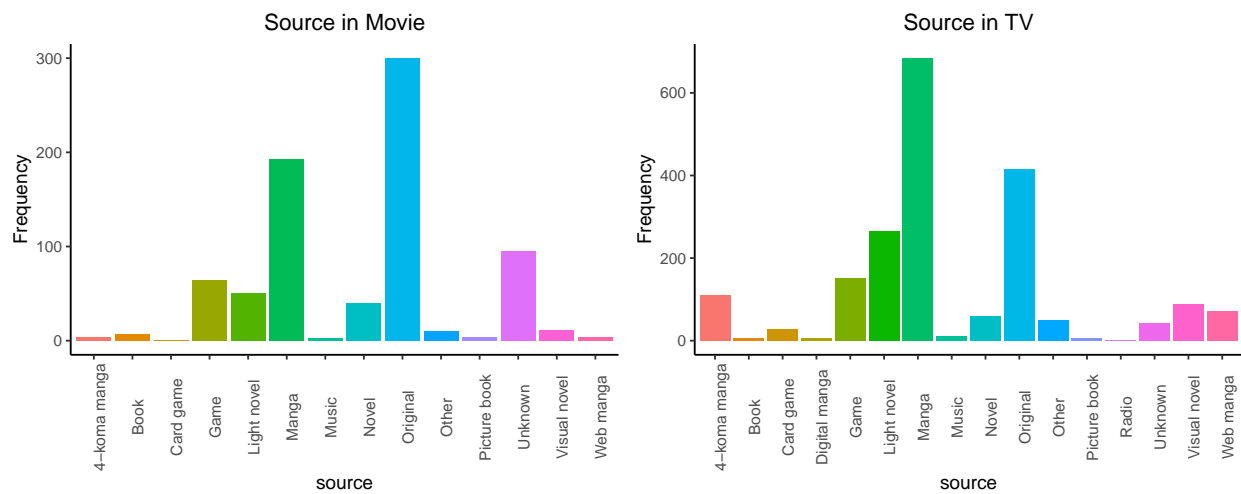


Figure 2: Frequency of Anime in Different Sources

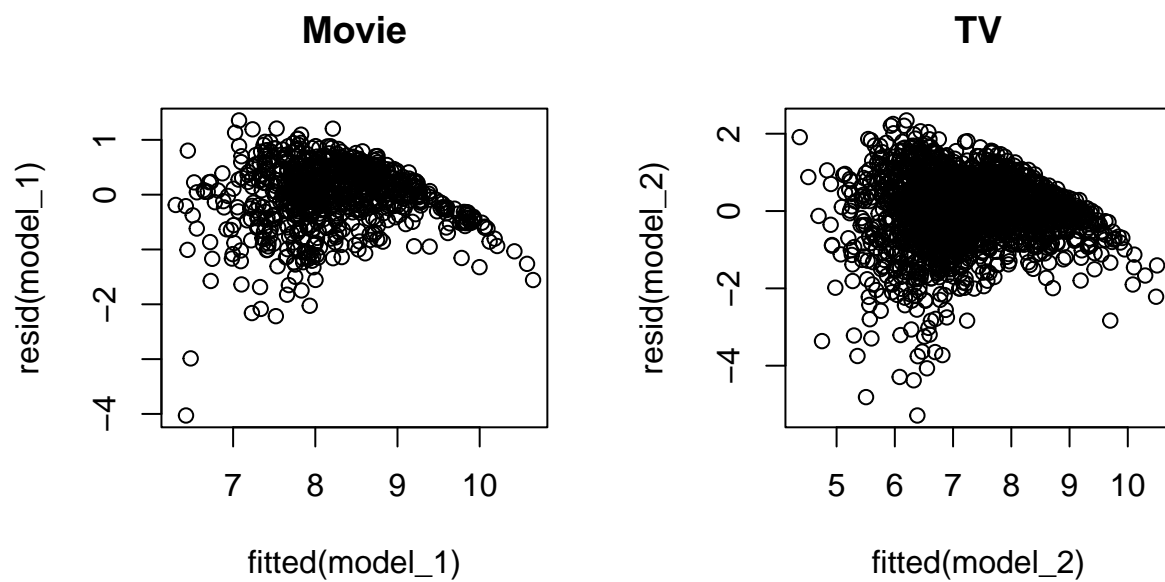


Figure 3: Residual Plot

For two two different subsets I fitted two different multilevel models, the variables I use for multilevel are rating, score, rank, and *rating : score*, for the random effect is ($1|rating : episodes$). I also used two individual plots to check the model, “Residuals vs. Fitted” and “QQ Plot”. The two multilevel models can be written as follow:

$$\begin{aligned} \log(popularity) = & 15.95 + 1.198 \cdot ratingNone + 0.979 \cdot ratingPG + 2.476 \cdot ratingPG13 + 1.215 \cdot ratingR \\ & - 0.935 \cdot ratingR+ - 0.97 \cdot score - 0.0001 \cdot rank + 0.195 \cdot ratingNone : score \\ & + 0.145 \cdot ratingPG : score - 0.397 \cdot ratingPG - 13 : score + 0.219 \cdot ratingR : score \\ & + 0.017 \cdot ratingR+ : score + n_j + \epsilon \\ n_j \sim & N(0, \sigma_a^2) \end{aligned}$$

where n_j is random effect: *rating : episodes*

$$\begin{aligned} \log(popularity) = & 12.67 - 2.891 \cdot ratingNone - 1.796 \cdot ratingPG + 1.022 \cdot ratingPG13 + 1.077 \cdot ratingR \\ & - 0.317 \cdot ratingR+ - 0.614 \cdot score - 0.0001 \cdot rank + 0.382 \cdot ratingNone : score \\ & + 0.274 \cdot ratingPG : score - 0.263 \cdot ratingPG - 13 : score + 0.369 \cdot ratingR : score \\ & + 0.174 \cdot ratingR+ : score + n_j + \epsilon \\ n_j \sim & N(0, \sigma_a^2) \end{aligned}$$

where n_j is random effect: *rating : episodes*

Therefore, we can interpret these two models by using, for instance, rating levels. For every one unit change in ratingR+, when other variables are constant, the $\log(popularity)$ will increase -0.169 in Movie. Same interpretation with TV, for every one unit change in ratingR+, when other variables are constant, the $\log(popularity)$ will increase in -1.512.

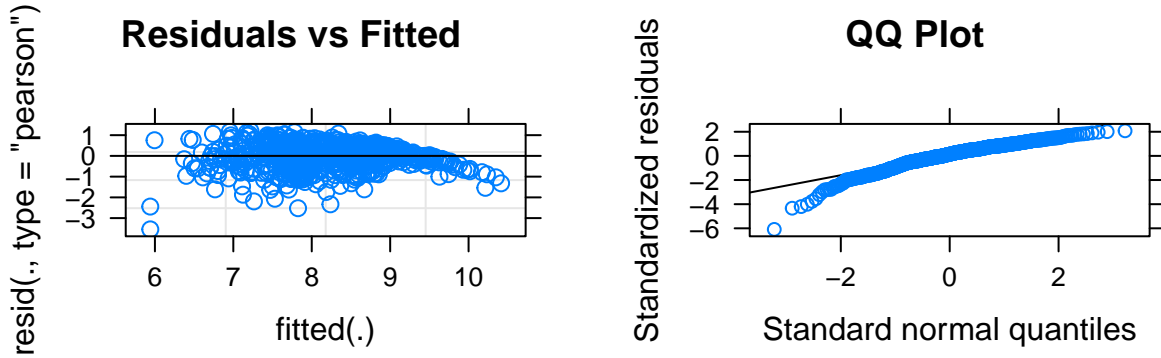


Figure 4: Type of Anime: Movie

According to the residual analysis, even though I used two different subsets to fit the multilevel model, the “Residuals vs. Fitted” plots indicated that the residual is not good to fit the model, so it might conclude there is not a lot of correlation between the predictors. On the QQ plots, there are plenty of residual points not on the lines, so it might not follow the normal distribution.

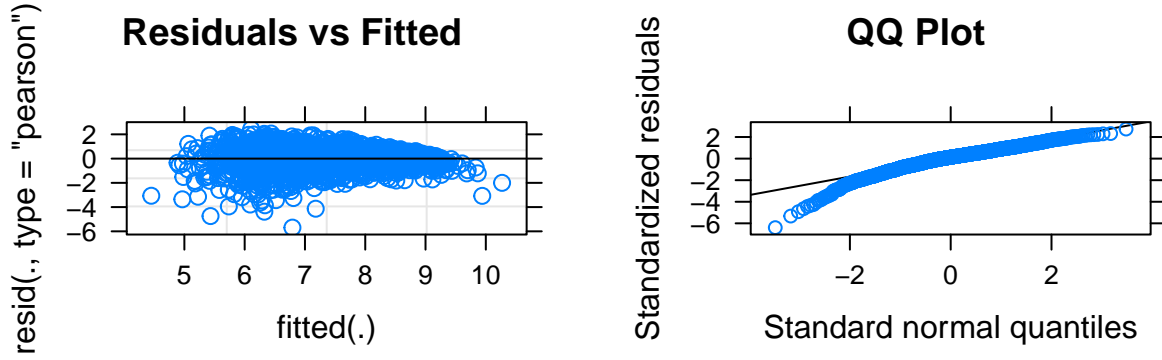


Figure 5: Type of Anime: TV

Result

Focusing on the model fitting and EDA I have done above, unfortunately, it is hard to identify the variables that might affect the popularity of Anime. In other words, these variables do not have a significant effect on the final results. Although I derived the Anime genre share from 2008 to 2020 from different subsets in figure 2, I still could not reach a valid conclusion. Moreover, I have analyzed rating and even though it has a significant effect in all variables, it has little effect on the overall model. The current result I have for two individual model is: 1) Movie: for every one unit change in ratingR+, when other variables are constant, the $\log(\text{popularity})$ will increase -0.169 in Movie; 2) TV: for every one unit change in ratingR+, when other variables are constant, the $\log(\text{popularity})$ will increase in -1.512.

Discussion

Overall the entire analysis was not good enough, even though the data looks good. The data I selected from Kaggle, I used rating, scour, rank as variables and source:genres as random effect for my multilevel model. Unfortunately, as I mentioned above it is unlikely to answer the questions I ask at the beginning of the report. After checking the model and data, I think there are some limitations for the data. First, the data provide the source level, which is helpful, I thought, for the model fitting. However, it did not provide the feedback of the user to indicate why it is so popular. Second, the variables I used might not be enough for model fitting, but the rest of variables were hard to fill in the lmer function. Last but not least, even though we have the count for type of Anime, we do not have actual market sale data for the Anime, including the type, rating, genres. This data might help to fit a better model in the future.

For the next step, I would like to collect the actual market data that I just mentioned, also I will try to include the market share data since it is more reasonable to indicate the popularity of Anime. In addition, I would like to learn more about modeling fitting and model design, because I realized I have lacked experience with designing a good model for data analysis.

Appendix

More EDA

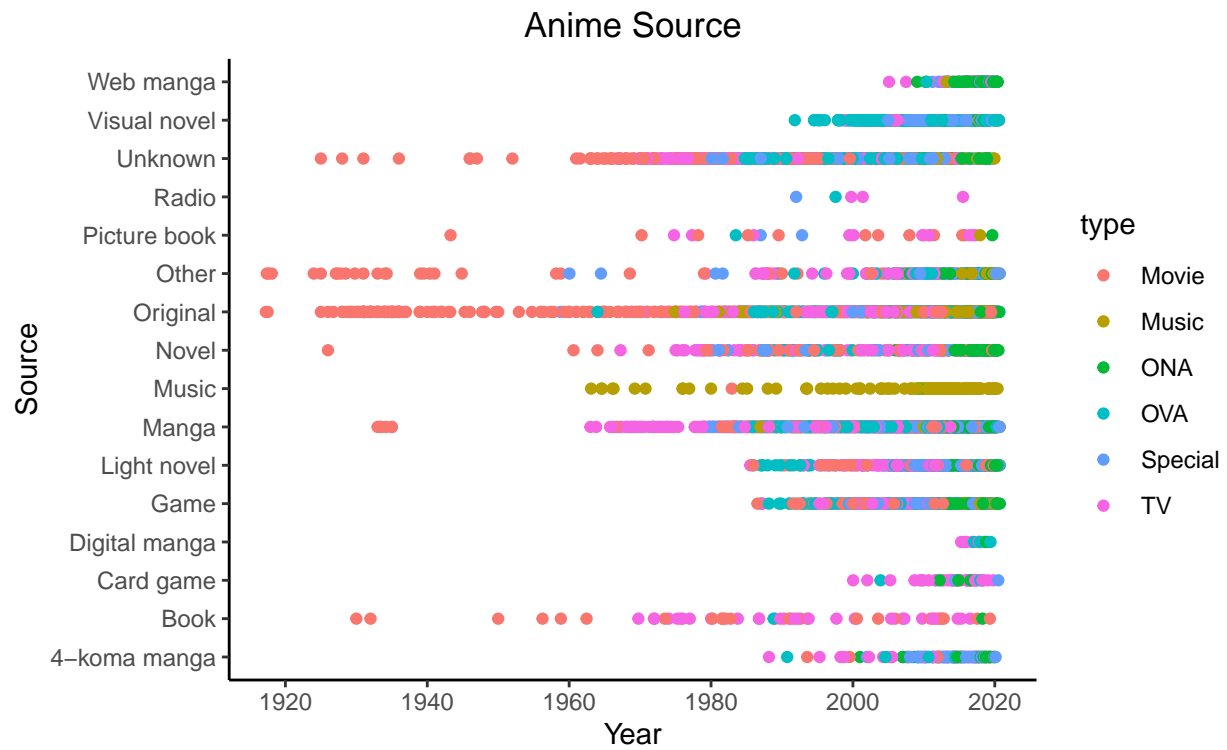


Figure 6: Original Dataset

```
## stan_glm
## family:      gaussian [identity]
## formula:     log(popularity) ~ rating + score + source + rank
## observations: 785
## predictors:  21
## -----
##
##               Median MAD_SD
## (Intercept)      16.5    1.0
## ratingNone         0.0    0.2
## ratingPG - Children  0.0    0.1
## ratingPG-13 - Teens 13 or older -0.4    0.1
## ratingR - 17+ (violence & profanity) -0.4    0.1
## ratingR+ - Mild Nudity -0.6    0.1
## score             -1.1    0.1
## sourceBook         0.2    0.4
## sourceCard game    0.2    0.7
## sourceGame        -0.1    0.3
## sourceLight novel  -0.4    0.3
## sourceManga       -0.1    0.3
## sourceMusic       -0.6    0.4
## sourceNovel        0.0    0.3
## sourceOriginal    -0.1    0.3
```

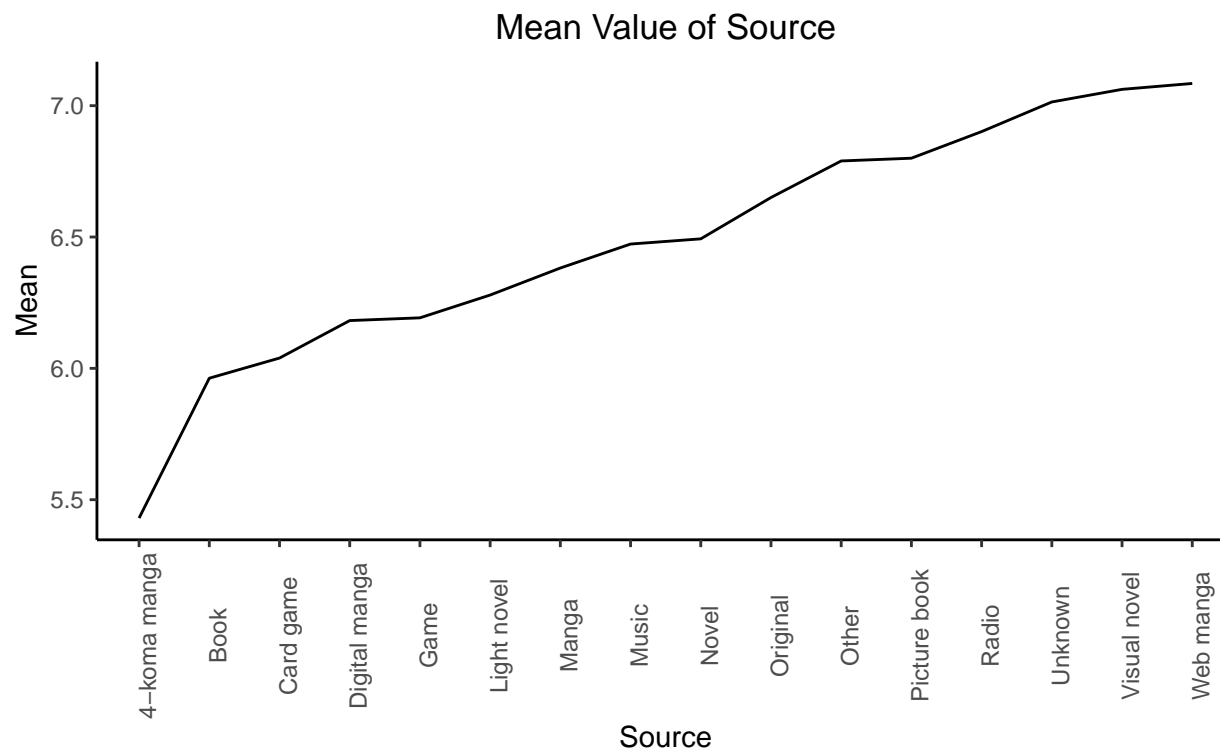


Figure 7: Mean Value of Each Source

```
## sourceOther                -0.3    0.4
## sourcePicture book         0.4    0.4
## sourceUnknown              0.1    0.3
## sourceVisual novel        -0.3    0.4
## sourceWeb manga           0.0    0.4
## rank                       0.0    0.0
##
## Auxiliary parameter(s):
##   Median MAD_SD
## sigma 0.6    0.0
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg

## stan_glm
## family:      gaussian [identity]
## formula:     log(popularity) ~ rating + score + source + rank
## observations: 1996
## predictors:  23
## -----
##
##               Median MAD_SD
## (Intercept)   15.4    0.9
## ratingNone    -0.5    0.3
## ratingPG - Children    0.0    0.1
## ratingPG-13 - Teens 13 or older -0.8    0.1
```

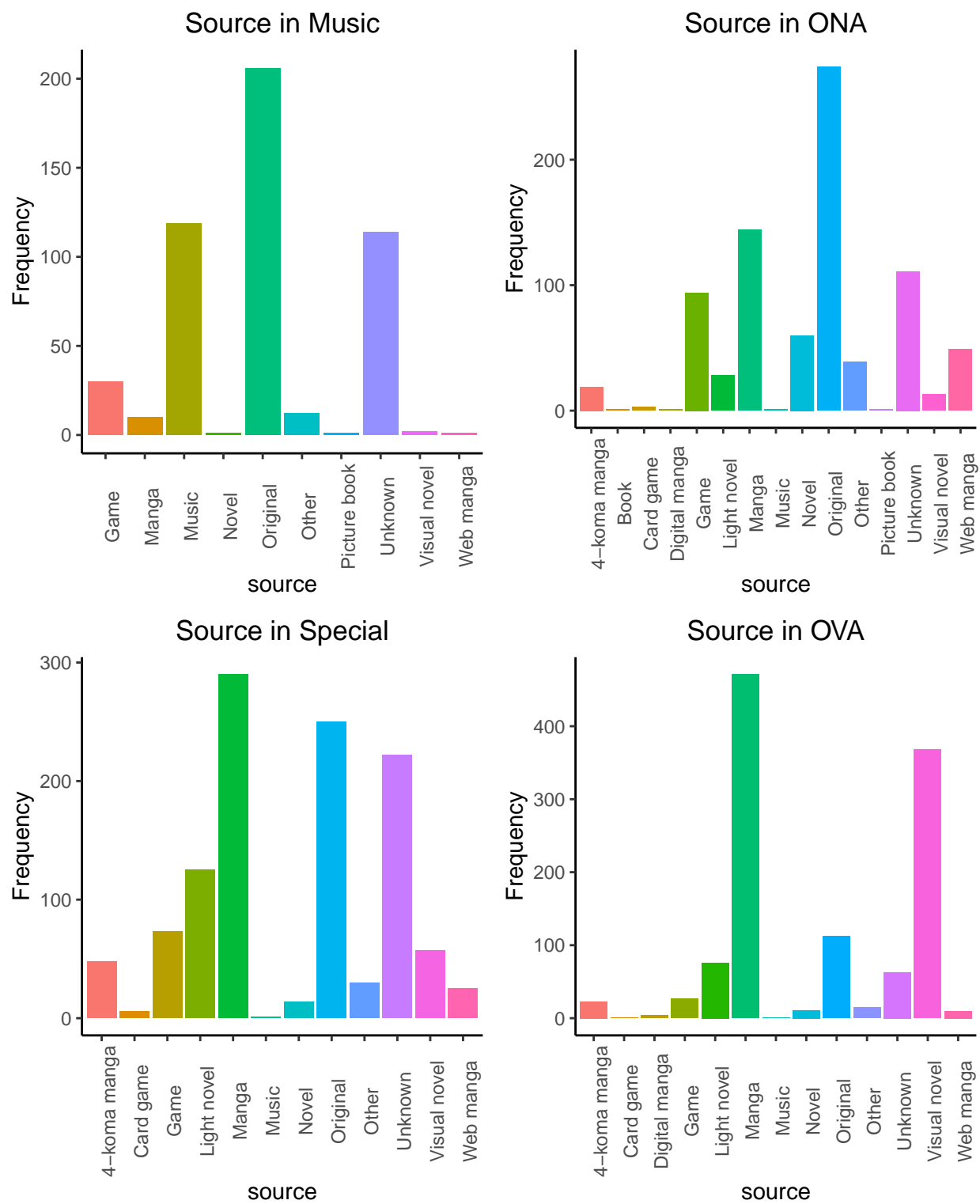


Figure 8: Different Type of Anime

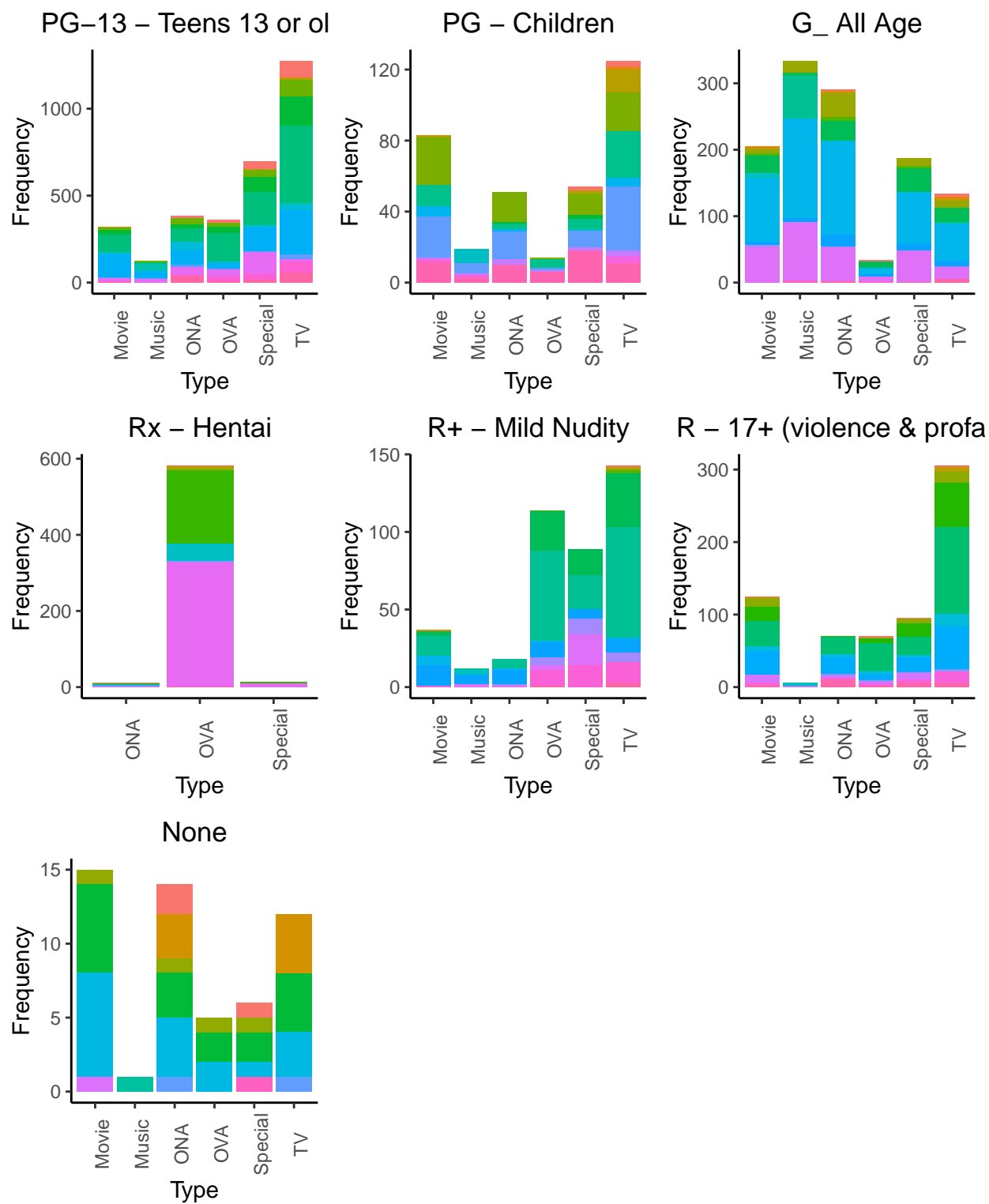


Figure 9: Rating Level

```

## ratingR - 17+ (violence & profanity) -1.3    0.1
## ratingR+ - Mild Nudity -1.2    0.1
## score -1.0    0.1
## sourceBook 0.3    0.4
## sourceCard game 0.4    0.2
## sourceDigital manga -0.5    0.3
## sourceGame 0.0    0.1
## sourceLight novel -0.8    0.1
## sourceManga -0.1    0.1
## sourceMusic 0.2    0.3
## sourceNovel 0.4    0.1
## sourceOriginal 0.1    0.1
## sourceOther 0.1    0.1
## sourcePicture book -0.1    0.4
## sourceRadio 0.3    0.9
## sourceUnknown 0.4    0.2
## sourceVisual novel -0.3    0.1
## sourceWeb manga -0.1    0.1
## rank 0.0    0.0
##
## Auxiliary parameter(s):
##      Median MAD_SD
## sigma 0.9    0.0
##
## -----
## * For help interpreting the printed output see ?print.stanreg
## * For info on the priors used see ?prior_summary.stanreg

## Linear mixed model fit by REML ['lmerMod']
## Formula: log(popularity) ~ rating + score + rank + rating:score + (1 |
##      rating:episodes)
##      Data: movie
## REML criterion at convergence: 1429.893
## Random effects:
##      Groups      Name      Std.Dev.
##      rating:episodes (Intercept) 0.1884
##      Residual 0.5819
## Number of obs: 785, groups: rating:episodes, 22
## Fixed Effects:
##              (Intercept)
##              15.9500733
##              ratingNone
##              1.1980332
##              ratingPG - Children
##              0.9798849
##              ratingPG-13 - Teens 13 or older
##              2.4762394
##              ratingR - 17+ (violence & profanity)
##              1.2146646
##              ratingR+ - Mild Nudity
##              -0.9345315
##              score
##              -0.9700600
##              rank

```

```

##                -0.0001564
##                ratingNone:score
##                -0.1946988
##                ratingPG - Children:score
##                -0.1453592
##                ratingPG-13 - Teens 13 or older:score
##                -0.3970980
## ratingR - 17+ (violence & profanity):score
##                -0.2191430
##                ratingR+ - Mild Nudity:score
##                0.0172206
## fit warnings:
## Some predictor variables are on very different scales: consider rescaling

## Linear mixed model fit by REML ['lmerMod']
## Formula: log(popularity) ~ rating + score + rank + rating:score + (1 |
##          rating:episodes)
## Data: tv
## REML criterion at convergence: 5288.396
## Random effects:
## Groups      Name      Std.Dev.
## rating:episodes (Intercept) 0.2468
## Residual      0.8900
## Number of obs: 1996, groups: rating:episodes, 154
## Fixed Effects:
##                (Intercept)
##                1.267e+01
##                ratingNone
##                -2.891e+00
##                ratingPG - Children
##                -1.796e+00
##                ratingPG-13 - Teens 13 or older
##                1.022e+00
##                ratingR - 17+ (violence & profanity)
##                1.077e+00
##                ratingR+ - Mild Nudity
##                -3.165e-01
##                score
##                -6.140e-01
##                rank
##                -7.334e-06
##                ratingNone:score
##                3.816e-01
##                ratingPG - Children:score
##                2.742e-01
##                ratingPG-13 - Teens 13 or older:score
##                -2.630e-01
## ratingR - 17+ (violence & profanity):score
##                -3.688e-01
##                ratingR+ - Mild Nudity:score
##                -1.744e-01
## fit warnings:
## Some predictor variables are on very different scales: consider rescaling

```

Random Effect in Two Models:

```
## $`rating:episodes`
##                               (Intercept)
## G - All Ages:1                5.461915e-02
## G - All Ages:2               -4.409337e-02
## G - All Ages:3                1.605115e-02
## G - All Ages:4               -8.353025e-03
## G - All Ages:8               -1.275184e-02
## G - All Ages:12              -5.472066e-03
## None:1                       1.541590e-02
## None:4                       -1.541590e-02
## PG - Children:1              2.201067e-13
## PG-13 - Teens 13 or older:1 -1.655047e-01
## PG-13 - Teens 13 or older:2  1.010908e-01
## PG-13 - Teens 13 or older:3  1.897905e-02
## PG-13 - Teens 13 or older:4  7.301148e-02
## PG-13 - Teens 13 or older:6 -4.539364e-02
## PG-13 - Teens 13 or older:10 1.781703e-02
## R - 17+ (violence & profanity):1 -2.018609e-01
## R - 17+ (violence & profanity):3 1.186265e-02
## R - 17+ (violence & profanity):4 1.978997e-01
## R - 17+ (violence & profanity):5 1.015224e-02
## R - 17+ (violence & profanity):6 -1.805369e-02
## R+ - Mild Nudity:1           1.771175e-02
## R+ - Mild Nudity:2          -1.771175e-02
##
## with conditional variances for "rating:episodes"
```

```
## $`rating:episodes`
##                               (Intercept)
## G - All Ages:3               -0.008412616
## G - All Ages:4               -0.049207920
## G - All Ages:5               -0.003327991
## G - All Ages:9               -0.028038682
## G - All Ages:10              -0.022028886
## G - All Ages:11              0.019637005
## G - All Ages:12              -0.133704899
## G - All Ages:13              -0.073557797
## G - All Ages:15              0.004230827
## G - All Ages:20              0.015046182
## G - All Ages:21              -0.094023779
## G - All Ages:22              0.013038449
## G - All Ages:24              0.020437533
## G - All Ages:25              -0.118855998
## G - All Ages:26              0.014443968
## G - All Ages:29              0.015671991
## G - All Ages:30              -0.002464677
## G - All Ages:32              0.054599344
## G - All Ages:34              0.031066695
## G - All Ages:36              0.035469630
## G - All Ages:37              0.041438014
## G - All Ages:39              -0.006446964
## G - All Ages:43              -0.038930461
```

## G - All Ages:44	0.021091881
## G - All Ages:46	-0.061156126
## G - All Ages:47	-0.047196725
## G - All Ages:48	0.009886584
## G - All Ages:49	0.073821066
## G - All Ages:50	0.077438559
## G - All Ages:51	0.036366272
## G - All Ages:52	0.112626664
## G - All Ages:58	0.035811720
## G - All Ages:59	0.036138863
## G - All Ages:60	0.021818964
## G - All Ages:78	0.024568444
## G - All Ages:97	-0.027295130
## None:10	0.031150746
## None:12	-0.071237320
## None:13	0.048724006
## None:24	-0.045506018
## None:26	0.036868586
## PG - Children:4	-0.011192590
## PG - Children:8	0.044304505
## PG - Children:10	0.015611008
## PG - Children:11	0.024736011
## PG - Children:12	-0.051153725
## PG - Children:13	-0.035956700
## PG - Children:14	-0.041644515
## PG - Children:15	-0.022580144
## PG - Children:16	-0.025244313
## PG - Children:20	-0.010241843
## PG - Children:21	0.044163798
## PG - Children:22	-0.027740467
## PG - Children:23	-0.004171472
## PG - Children:24	-0.077010857
## PG - Children:25	-0.056740459
## PG - Children:26	0.016940128
## PG - Children:29	0.056474670
## PG - Children:30	-0.050361140
## PG - Children:36	0.011027156
## PG - Children:37	0.005899504
## PG - Children:38	0.004649234
## PG - Children:39	0.043406699
## PG - Children:47	-0.019875420
## PG - Children:48	0.024955667
## PG - Children:49	0.037804252
## PG - Children:50	0.122896985
## PG - Children:51	0.045402947
## PG - Children:52	0.077290073
## PG - Children:60	-0.005016620
## PG - Children:64	0.004002441
## PG - Children:75	-0.012178186
## PG - Children:76	0.032226690
## PG - Children:84	-0.083429736
## PG - Children:93	-0.079638886
## PG - Children:100	0.002385307
## PG-13 - Teens 13 or older:3	0.054491120

## PG-13 - Teens 13 or older:4	0.030049731
## PG-13 - Teens 13 or older:5	0.052174907
## PG-13 - Teens 13 or older:6	0.268634297
## PG-13 - Teens 13 or older:8	-0.031642848
## PG-13 - Teens 13 or older:9	0.091341208
## PG-13 - Teens 13 or older:10	-0.317452621
## PG-13 - Teens 13 or older:11	-0.055587026
## PG-13 - Teens 13 or older:12	-0.363982925
## PG-13 - Teens 13 or older:13	-0.181724224
## PG-13 - Teens 13 or older:14	0.095832574
## PG-13 - Teens 13 or older:16	-0.009993254
## PG-13 - Teens 13 or older:17	0.075809615
## PG-13 - Teens 13 or older:18	0.109089404
## PG-13 - Teens 13 or older:20	-0.248305480
## PG-13 - Teens 13 or older:21	0.030778773
## PG-13 - Teens 13 or older:22	-0.255893068
## PG-13 - Teens 13 or older:23	0.004652431
## PG-13 - Teens 13 or older:24	-0.459917728
## PG-13 - Teens 13 or older:25	-0.473704931
## PG-13 - Teens 13 or older:26	0.118219102
## PG-13 - Teens 13 or older:30	0.061763133
## PG-13 - Teens 13 or older:31	0.083909016
## PG-13 - Teens 13 or older:33	0.052225032
## PG-13 - Teens 13 or older:34	0.122332454
## PG-13 - Teens 13 or older:36	-0.001508934
## PG-13 - Teens 13 or older:37	0.224208312
## PG-13 - Teens 13 or older:38	0.044613497
## PG-13 - Teens 13 or older:39	0.153417158
## PG-13 - Teens 13 or older:40	0.054931988
## PG-13 - Teens 13 or older:48	0.036873936
## PG-13 - Teens 13 or older:49	0.061934951
## PG-13 - Teens 13 or older:50	0.072714815
## PG-13 - Teens 13 or older:51	0.057492681
## PG-13 - Teens 13 or older:52	0.303475485
## PG-13 - Teens 13 or older:59	0.064461720
## PG-13 - Teens 13 or older:60	-0.085159100
## PG-13 - Teens 13 or older:61	0.010450186
## PG-13 - Teens 13 or older:65	0.022360241
## PG-13 - Teens 13 or older:73	-0.040727782
## PG-13 - Teens 13 or older:75	0.017970831
## PG-13 - Teens 13 or older:77	0.012697978
## PG-13 - Teens 13 or older:89	0.047829817
## PG-13 - Teens 13 or older:97	-0.060154869
## PG-13 - Teens 13 or older:99	0.046839646
## PG-13 - Teens 13 or older:100	0.102178752
## R - 17+ (violence & profanity):3	0.037799670
## R - 17+ (violence & profanity):4	0.008765925
## R - 17+ (violence & profanity):5	0.010443608
## R - 17+ (violence & profanity):6	0.032683310
## R - 17+ (violence & profanity):7	0.093285866
## R - 17+ (violence & profanity):8	0.035934270
## R - 17+ (violence & profanity):10	0.052797007
## R - 17+ (violence & profanity):11	-0.066414049
## R - 17+ (violence & profanity):12	-0.319511314

```

## R - 17+ (violence & profanity):13 0.125327138
## R - 17+ (violence & profanity):15 -0.106211933
## R - 17+ (violence & profanity):16 -0.074065268
## R - 17+ (violence & profanity):18 0.003405577
## R - 17+ (violence & profanity):22 -0.084500484
## R - 17+ (violence & profanity):23 0.032938672
## R - 17+ (violence & profanity):24 -0.129102234
## R - 17+ (violence & profanity):25 -0.008790274
## R - 17+ (violence & profanity):26 0.178124504
## R - 17+ (violence & profanity):27 0.072417101
## R - 17+ (violence & profanity):38 0.083387484
## R - 17+ (violence & profanity):39 0.256909549
## R - 17+ (violence & profanity):64 -0.235624123
## R+ - Mild Nudity:6 0.018343119
## R+ - Mild Nudity:8 0.204320291
## R+ - Mild Nudity:9 0.047641084
## R+ - Mild Nudity:10 -0.219686016
## R+ - Mild Nudity:11 0.181286546
## R+ - Mild Nudity:12 -0.195448437
## R+ - Mild Nudity:13 0.048539183
## R+ - Mild Nudity:24 -0.099856416
## R+ - Mild Nudity:26 -0.093141074
## R+ - Mild Nudity:50 0.108001722
##
## with conditional variances for "rating:episodes"

```