

Mini Project 2

Albert Chui (albertchui)

2024-05-08

Table of contents

1 Abstract	1
2 Dataset Examination	2
3 Testing for Differences in Mean Rating Across Genres	3
3.1 Documentary VS Action	3
3.2 Crime vs. Sports	4
3.3 News vs Drama	5
3.4 Romance vs Sci-Fi vs War vs Western	6
3.4.1 Check Assumptions	6
3.5 Conduct ANOVA Test	9
3.6 Fantasy vs Comedy vs Horror	9
3.6.1 Check Assumptions	9
3.7 Conduct ANOVA Test	12
4 Testing for Differences in Mean Ratings Within Genres Across Years	12
4.1 Conduct ANOVA Test for Shorts, History, and Animation Mean Ratings Across Years	13
4.1.1 Check Assumptions	13
4.2 Conduct ANOVA Test	13

1 Abstract

In this paper we explore data from the official IMdB Documentation detailing 942007 movies. We find that the average rating between genres varies between 5 and 7 points. For the ratings of selected genres of news, reality-tv, horror, documentary we see that documentaries are relativly stable, horror has been decreasing, reality-tv has been increasing, and news has been extremly variable

throughout the years. In addition, runtime of movies has been increasing over the years, to a current 2024 average of 120 minutes. The runtime of episodes in tv series had been increasing until 24 years ago, when it has suddenly started decreasing until it hit 40 minutes, then bounced back to 50 minutes.

2 Dataset Examination

Let us first look at all of our unique genres in this dataset:

genres	avg_rating	num_movies
Documentary	7.235	163250
Short	7.126	4243
Biography	6.975	20475
History	6.927	20115
Sport	6.767	14821
Animation	6.754	25459
Music	6.726	23800
Family	6.503	34319
War	6.47	10493
Film-Noir	6.467	882
Reality-TV	6.385	24960
Drama	6.364	293034
News	6.351	10100
Talk-Show	6.269	22822
Musical	6.247	12137
Romance	6.238	60294
Crime	6.187	49216
Adventure	6.186	39866
Fantasy	6.176	21801
Comedy	6.175	173668
Game-Show	6.128	8983
Mystery	6.094	23625
Western	5.954	8729
Action	5.945	66914
Thriller	5.729	51903
Sci-Fi	5.671	19062
Adult	5.561	11519
Horror	5.134	44503

3 Testing for Differences in Mean Rating Across Genres

In this section we will perform several hypothesis tests on different genres to determine if there is a difference in the mean rating based on genre.

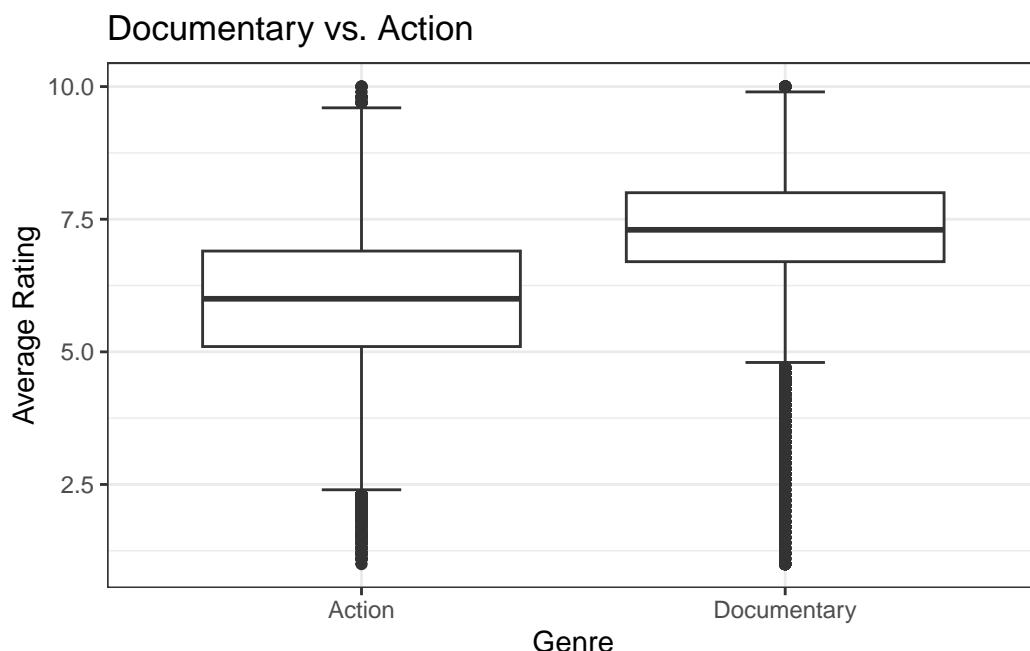
3.1 Documentary VS Action

Here, we specifically test for a difference in mean rating between documentaries and action programs. To do this we test the following hypotheses:

$$\begin{cases} H_0 : \mu_{documentary} = \mu_{action} \\ H_a : \mu_{documentary} \neq \mu_{action} \end{cases}$$

Welch Two Sample t-test

```
data: movies %>% filter(genres == "Documentary") %>% pull(averageRating) and movies %>%  
t = 149.26, df = 68285, p-value < 2.2e-16  
alternative hypothesis: true difference in means is not equal to 0  
95 percent confidence interval:  
 1.272481 1.306345  
sample estimates:  
mean of x mean of y  
 7.234642 5.945229
```



From the results of our two-sample t-test we can see that our p-value is smaller than 0.05, which means there is a significant difference in mean ratings across Documentaries and Action movies/series. From the boxplot, we can see that the mean rating for both categories seems to clearly differ.

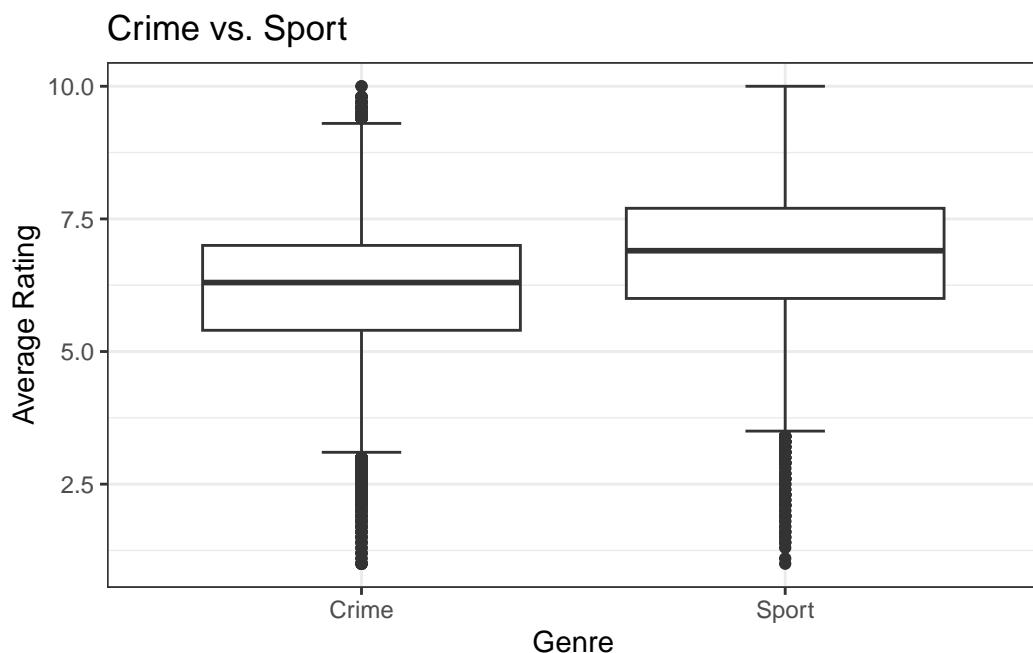
3.2 Crime vs. Sports

Here, we specifically test for a difference in mean rating between crime and sports programs. To do this we test the following hypotheses:

$$\begin{cases} H_0 : \mu_{\text{crime}} = \mu_{\text{sports}} \\ H_a : \mu_{\text{crime}} \neq \mu_{\text{sports}} \end{cases}$$

Welch Two Sample t-test

```
data: movies %>% filter(genres == "Crime") %>% pull(averageRating) and movies %>% filter(genres == "Sport") %>% pull(averageRating)
t = -31.313, df = 7909.9, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.6153849 -0.5428749
sample estimates:
mean of x mean of y
6.187413 6.766543
```



From the results of our two-sample t-test we can see that our p-value is smaller than 0.05, which means there is a significant difference in mean ratings across Crime and Sports movies/series. From the boxplot, we can see that the mean rating for both categories differs slightly, but only through our t-test can we tell that this is a significant difference.

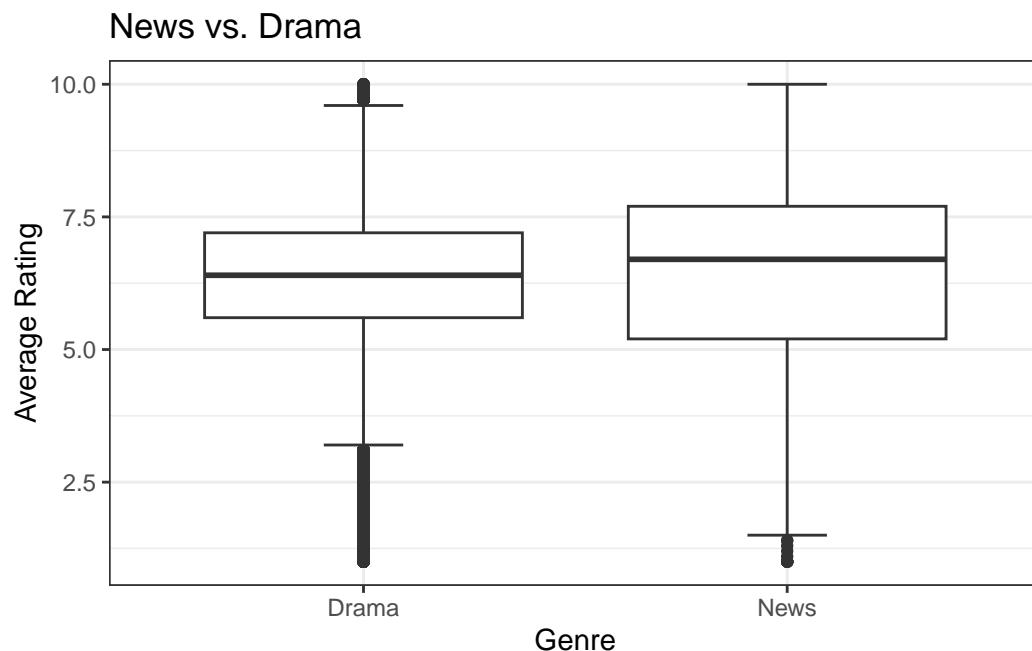
3.3 News vs Drama

Here, we specifically test for a difference in mean rating between News and Drama programs. To do this we test the following hypotheses:

$$\begin{cases} H_0 : \mu_{news} = \mu_{drama} \\ H_a : \mu_{news} \neq \mu_{drama} \end{cases}$$

Welch Two Sample t-test

```
data: movies %>% filter(genres == "News") %>% pull(averageRating) and movies %>% filter
t = -0.37348, df = 2709.3, p-value = 0.7088
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
-0.08026006 0.05457742
sample estimates:
mean of x mean of y
6.351294 6.364135
```



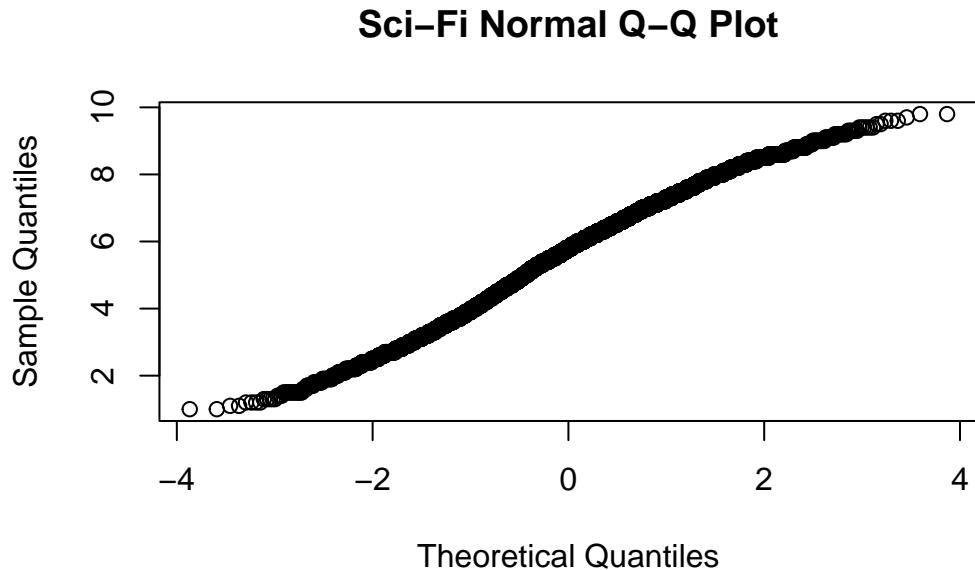
From the results of our two-sample t-test we can see that our p-value is larger than 0.05, which means there is not a significant difference in mean ratings across Crime and Sports movies/series. From the boxplot, we can see that the mean rating for both categories is very similar, appearing almost identical.

3.4 Romance vs Sci-Fi vs War vs Western

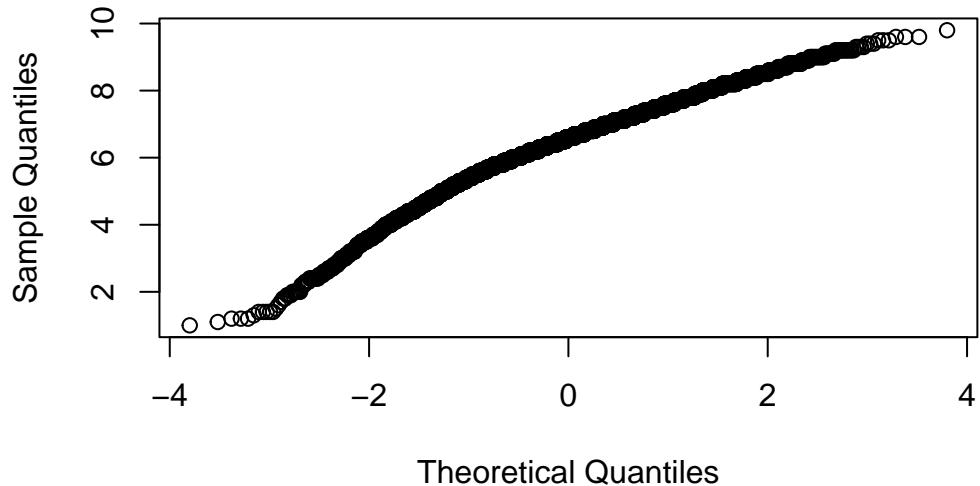
In this section we test for differences in mean ratings across groups larger than 2 genres, meaning we need to perform an Analysis of Variance test.

3.4.1 Check Assumptions

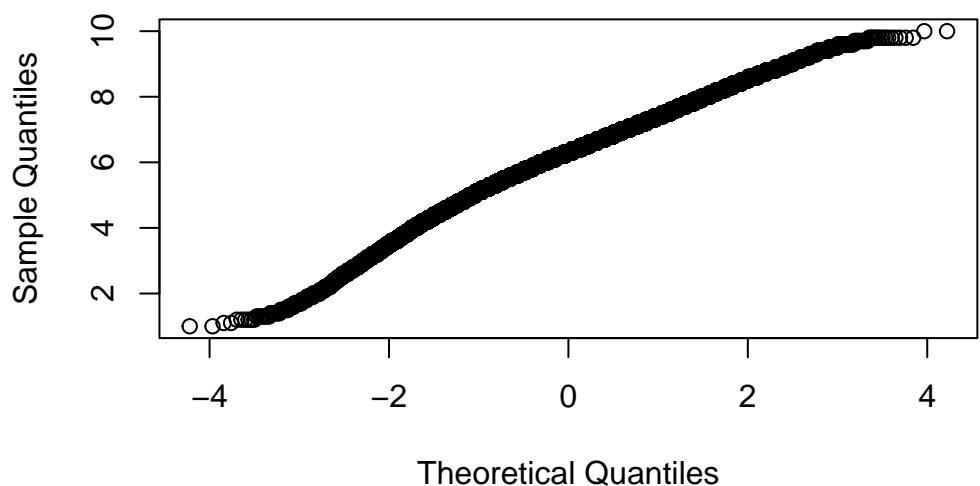
1) Normality



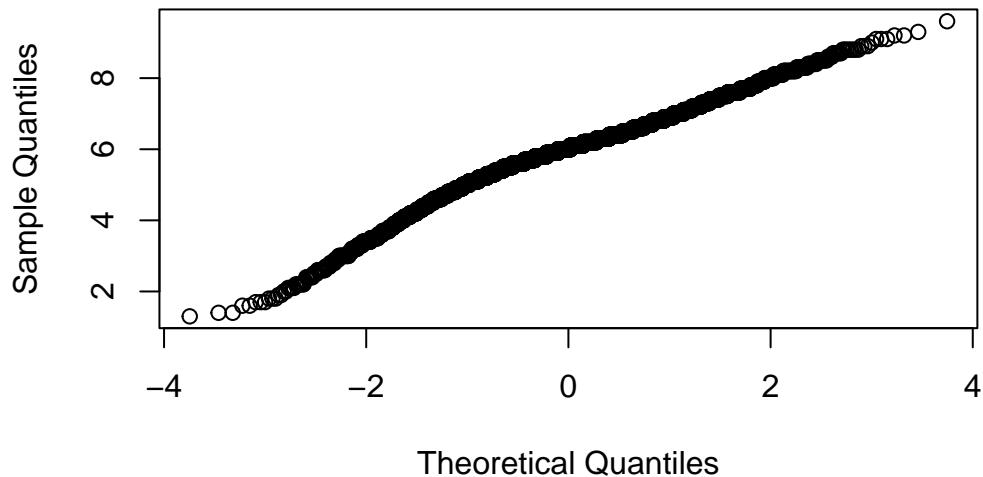
War Normal Q–Q Plot



Romance Normal Q–Q Plot

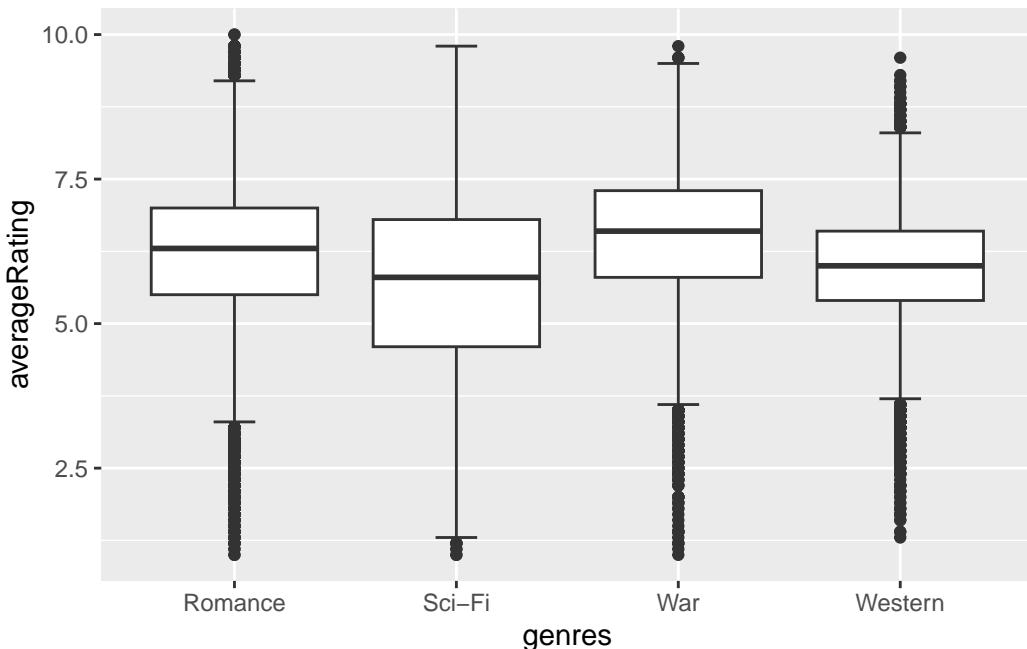


Western Normal Q-Q Plot



From the above Q-Q Plots, we can see that our normality assumption holds for all four genres.

2) Homoscedasticity (Constant Variance)



From the above boxplots, we can see that the assumption of Homoscedasticity holds for this dataset as all four genres seem to have approximately the same variance.

3) Independence

Although we cannot test for independence, it is unlikely that the ratings for one movie or series strongly effects the rating for a different movie or series, meaning our assumption of Independence is most likely not violated.

3.5 Conduct ANOVA Test

$$\begin{cases} H_0 : \mu_{war} = \mu_{romance} = \mu_{western} = \mu_{sci-fi} \\ H_a : \text{At least one mean differs from the others} \end{cases}$$

```
Df Sum Sq Mean Sq F value Pr(>F)
genres      3    3328   1109.2   696.8 <2e-16 ***
Residuals  62743   99875        1.6
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the results of our ANOVA test above, we can see that our p-value is very small, much smaller than 0.05, meaning we can reject $H_0 : \mu_{war} = \mu_{romance} = \mu_{western} = \mu_{sci-fi}$, indicating there is a significant difference in the mean ratings for at least one of the tested genres (Romance, War, Western, Sci-Fi).

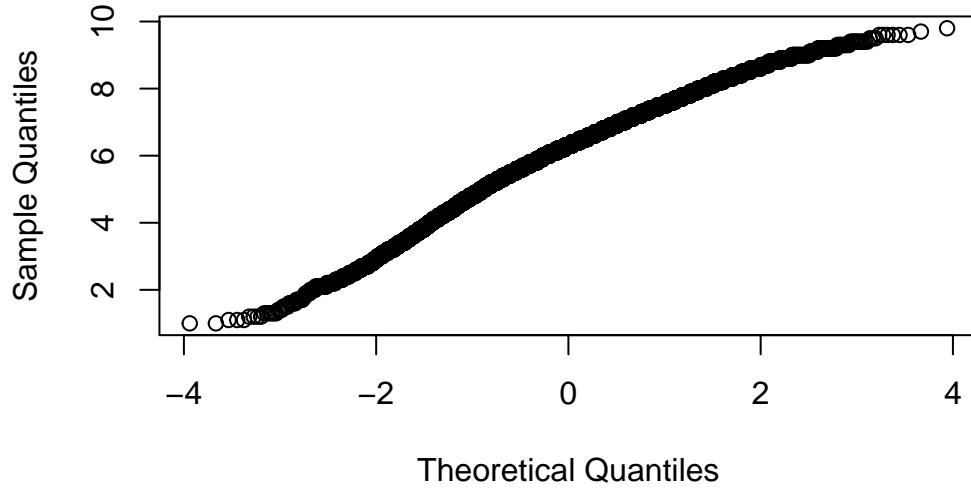
3.6 Fantasy vs Comedy vs Horror

In this section we test for differences in mean ratings across groups larger than 2 genres, meaning we need to perform an Analysis of Variance test.

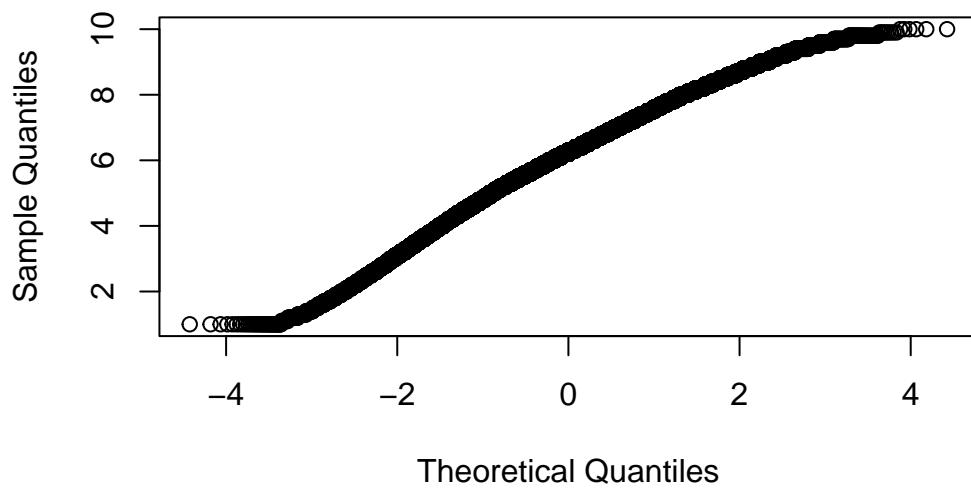
3.6.1 Check Assumptions

- 1) Normality

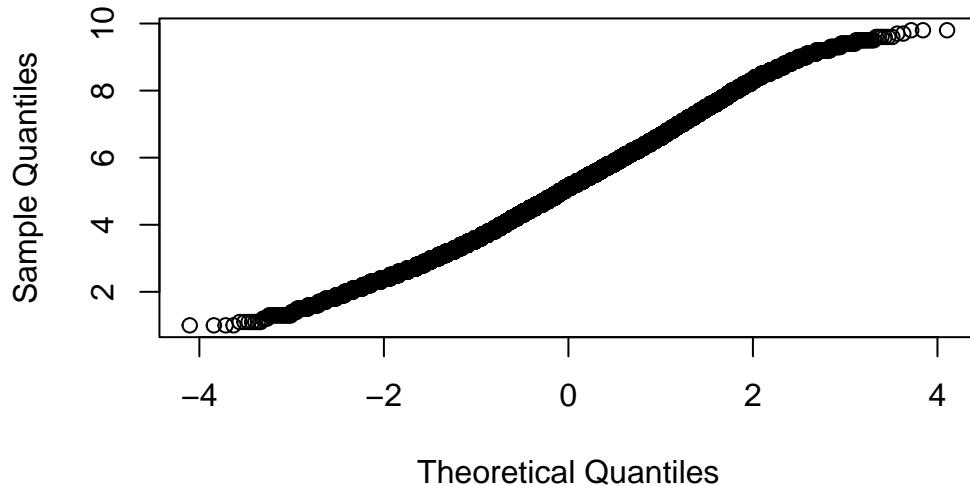
Fantasy Normal Q–Q Plot



Comedy Normal Q–Q Plot

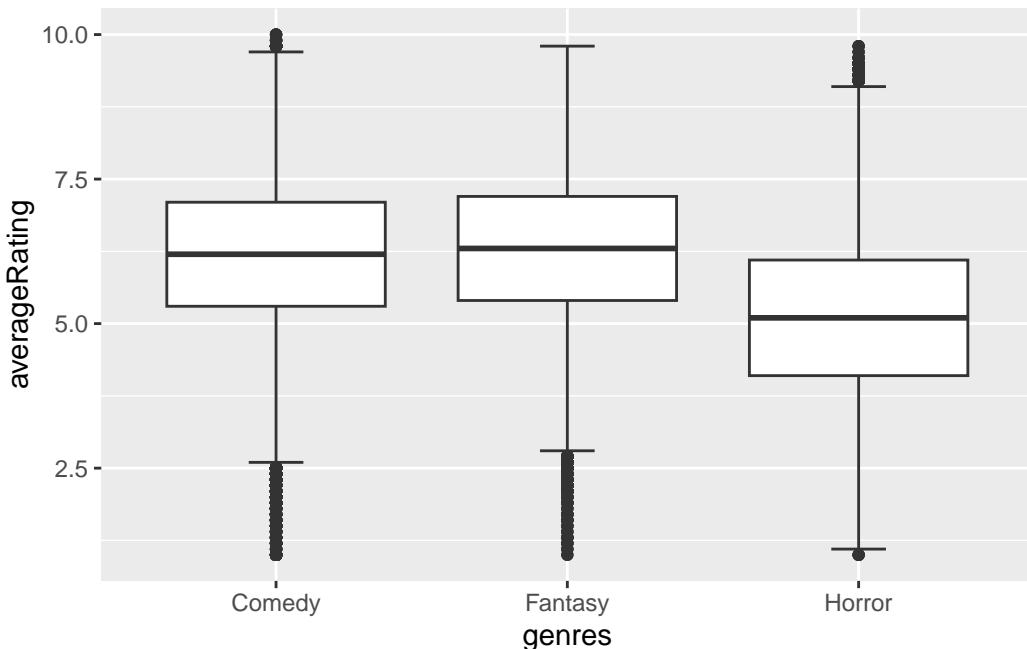


Horror Normal Q-Q Plot



From the above Q-Q Plots, we can see that our normality assumption holds for all three genres.

2) Homoscedasticity (Constant Variance)



From the above boxplots, we can see that the assumption of Homoscedasticity holds for this dataset as all three genres seem to have approximately the same variance.

3) Independence

Although we cannot test for independence, it is unlikely that the ratings for one movie or series strongly effects the rating for a different movie or series, meaning our assumption of Independence is most likely not violated.

3.7 Conduct ANOVA Test

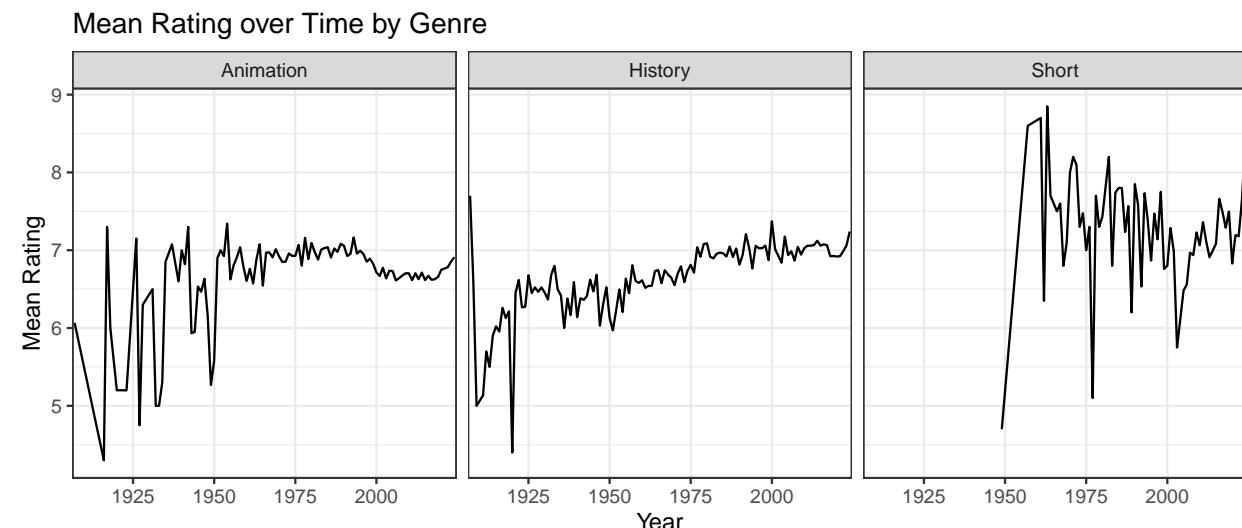
$$\begin{cases} H_0 : \mu_{fantasy} = \mu_{comedy} = \mu_{horror} \\ H_a : \text{At least one mean differs from the others} \end{cases}$$

```
Df Sum Sq Mean Sq F value Pr(>F)
genres      2  22058   11029     5566 <2e-16 ***
Residuals 140661 278717           2
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the results of our ANOVA test above, we can see that our p-value is very small, much smaller than 0.05, meaning we can reject $H_0 : \mu_{fantasy} = \mu_{comedy} = \mu_{horror}$, indicating there is a significant difference in the mean ratings for at least one of the tested genres (Fantasy, Comedy, Horror).

4 Testing for Differences in Mean Ratings Within Genres Across Years

In this section we investigate if there are significant changes to the mean ratings of several genres across different years. We focus on the genres of History, Short Films, and Animation.



From the above graph displaying mean ratings in each year by genre, we can already see that the mean rating within each genre seems to vary greatly each year and tends to oscillate especially before 1950. In order to statistically test for difference in mean rating across years, we conducted an ANOVA test for each genre.

4.1 Conduct ANOVA Test for Shorts, History, and Animation Mean Ratings Across Years

4.1.1 Check Assumptions

Since we have tested the assumption of Normality, Homoscedasticity, and Independence for other genres and found that they hold, it is reasonable to assume that they hold for these genres as well since all the data is coming from the same source.

4.2 Conduct ANOVA Test

$$\begin{cases} H_0 : \mu_i = \mu_j; \text{ for } i, j \in \text{Years and } i \neq j \\ H_a : \text{At least one mean differs from the others} \end{cases}$$

```
Df Sum Sq Mean Sq F value Pr(>F)
startYear     117    551   4.711   3.199 <2e-16 ***
Residuals  27061  39856   1.473
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

From the results of our ANOVA test above, we can see that our p-value is very small, much smaller than 0.05, meaning we can reject \$H_0: \{i\} = \{j\} ; \text{ for } i,j \in \text{Years and } i \neq j\$, indicating there is a significant difference in the mean ratings within genres across years for the tested genres (Shorts, Animation, and History).

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

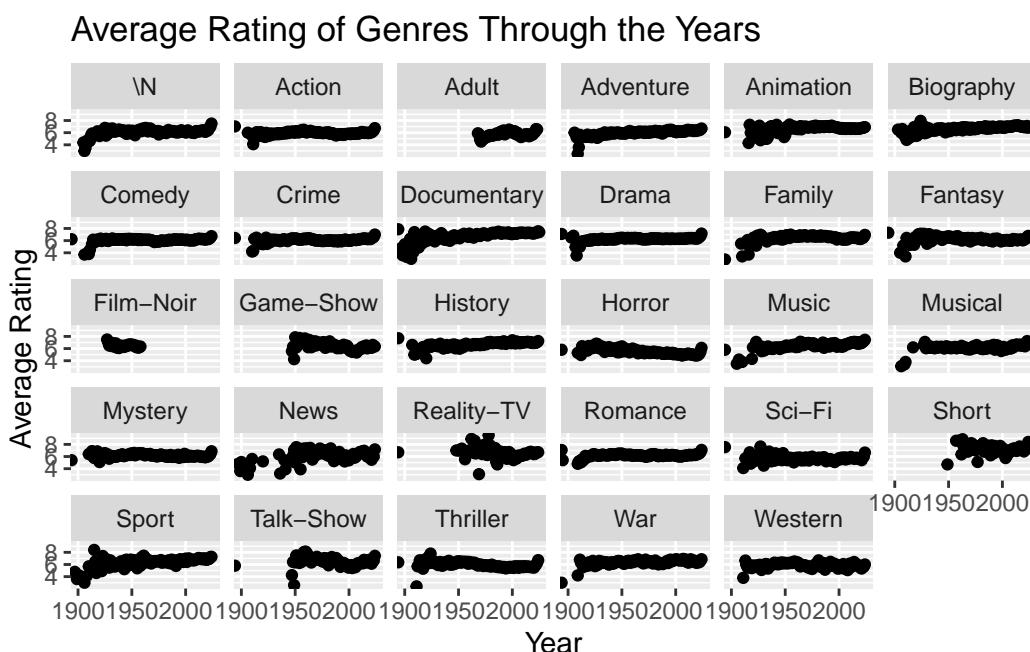
```
# A tibble: 3,146 x 3
# Groups:   genres [29]
  genres startYear avg_rating
  <chr>   <chr>        <dbl>
1 Action   1906         6
2 Action   1910         5.5
3 Action   1911         4.1
4 Action   1912         5.7
```

```

5 Action 1913      5.88
6 Action 1914      6.08
7 Action 1915      5.38
8 Action 1916      6.01
9 Action 1917      5.74
10 Action 1918     6.09
# i 3,136 more rows

```

Warning: Removed 151 rows containing missing values or values outside the scale range (`geom_point()`).



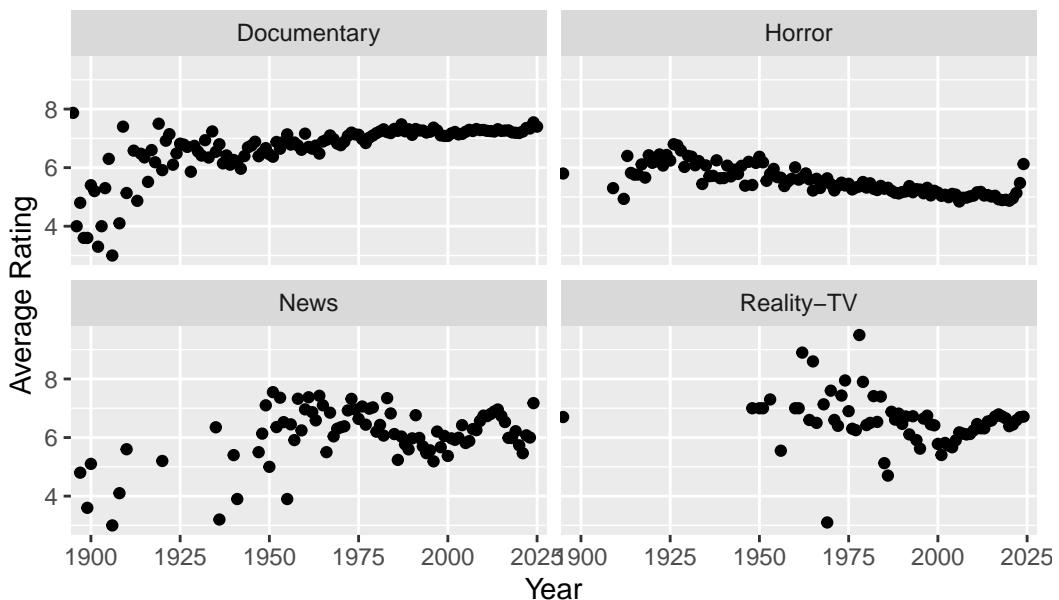
```

ggplot(limited_timed_genre_ratings, color = `genres`, aes(x = startYear, y = avg_rating,
  geom_point() +
  facet_wrap("genres") +
  labs(x = "Year", y = "Average Rating", title = "Average Rating of Genres Through the
  scale_x_discrete(breaks = seq(0, max(timed_genre_ratings$startYear), by = 25))

```

Warning: Removed 30 rows containing missing values or values outside the scale range (`geom_point()`).

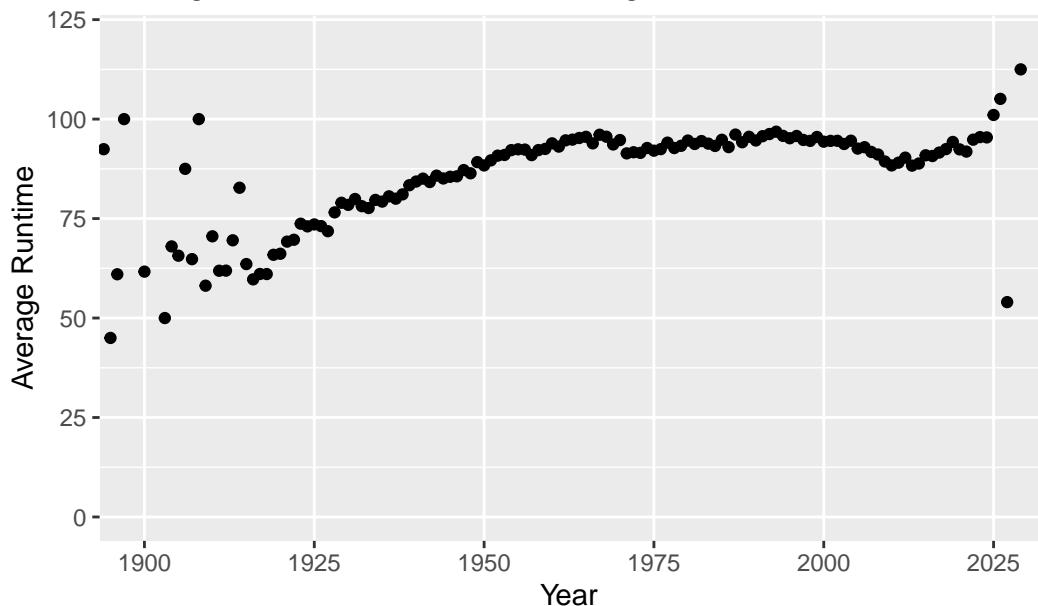
Average Rating of Genres Through the Years



Warning: NAs introduced by coercion

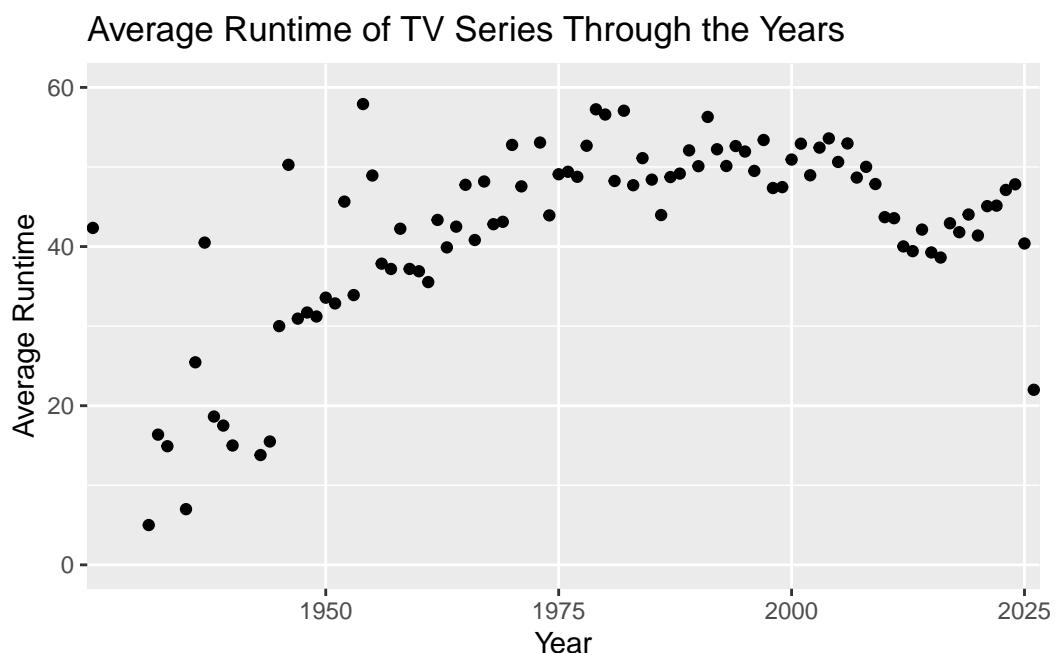
Warning: Removed 7 rows containing missing values or values outside the scale range (`geom_point()`).

Average Runtime of Movies Through the Years



```
ggplot(series_timed_runlength)+  
  geom_point(aes(x=`startYear`, y = `avg_runtime`)) +  
  labs(x = "Year", y = "Average Runtime", title = "Average Runtime of TV Series Through the Years") +  
  scale_x_discrete(breaks = seq(0, max(series_timed_runlength$startYear), by = 25)) +  
  ylim(0, 60)
```

Warning: Removed 9 rows containing missing values or values outside the scale range
(`geom_point()`).



Joining with `by = join_by(startYear)`

Warning: Removed 2 rows containing missing values or values outside the scale range
(`geom_line()`).

Warning: Removed 7 rows containing missing values or values outside the scale range
(`geom_line()`).

Average Rating of Adult Movies vs Non–Adult Movies Through tl

