# Data 550: Data Visualization I

Lecture 4b: Comparing Distributions in R

## Dr. Irene Vrbik

University of British Columbia Okanagan

https://github.com/ubco-mds-2022/Data-550

1

# Introduction

- Up until this point we have provided examples mostly in Altair with the understanding that ggplot has a similar counterpart.

- As Altair is relatively new, and ggplot2 is one of the most widely used and documented packages in R, it does have functionalities that Altair has yet to implement.

- One such example is violin plots.

https://github.com/ubco-mds-2022/Data-550

# Learning Outcomes

- Create density, box plots, and violin plots using ggplot

https://github.com/ubco-mds-2022/Data-550

3

# Data

Below is the reprocessed movies data frame (to see how it was processed see the accompanying ipynb)

```r
1  # the above is the cleaned version
2  library(rjson)
3  library(tidyverse)
4  movies <- fromJSON(file = 'data/lec-movies.json') %>%
5      as_tibble() %>%
6      unnest(-c(countries, genres))
7
8  head(movies)
```

| id | title |
|----|-------|
| <dbl> | <chr> |
| 12 | Finding Nemo |
| 22 | Pirates of the Caribbean: The Curse of the Black Pearl |

https://github.com/ubco-mds-2022/Data-550

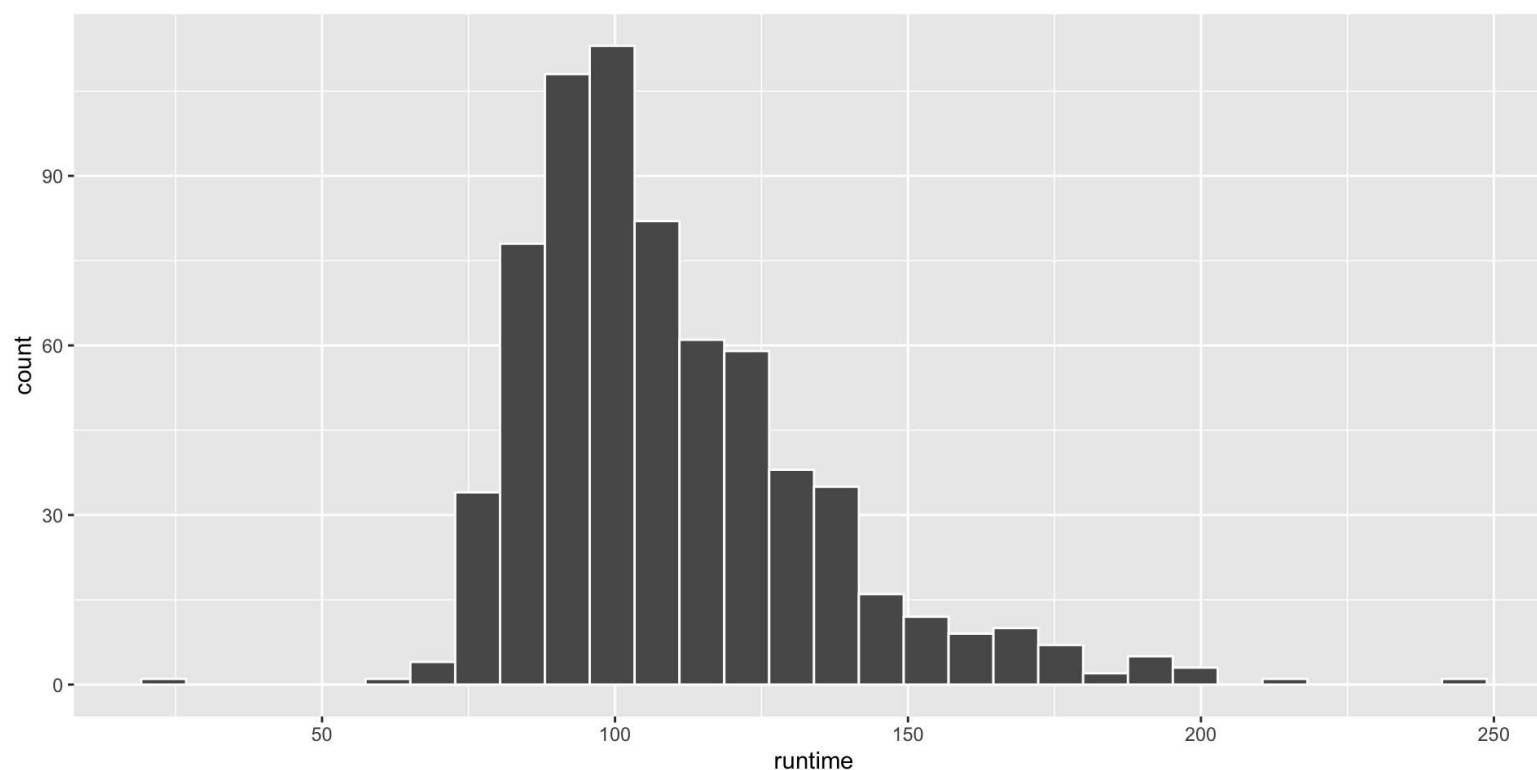| id | title |
| --- | --- |
| <dbl> | <chr> |
| 35 | The Simpsons Movie |
| 58 | Pirates of the Caribbean: Dead Man's Chest |
| 75 | Mars Attacks! |
| 117 | The Untouchables |

6 rows | 1-2 of 11 columns

https://github.com/ubco-mds-2022/Data-550

4

# Histogram

Let's recall how to make a histogram.

```
1  ggplot(movies, aes(x = runtime)) +
2      geom_histogram(color = 'white')
```



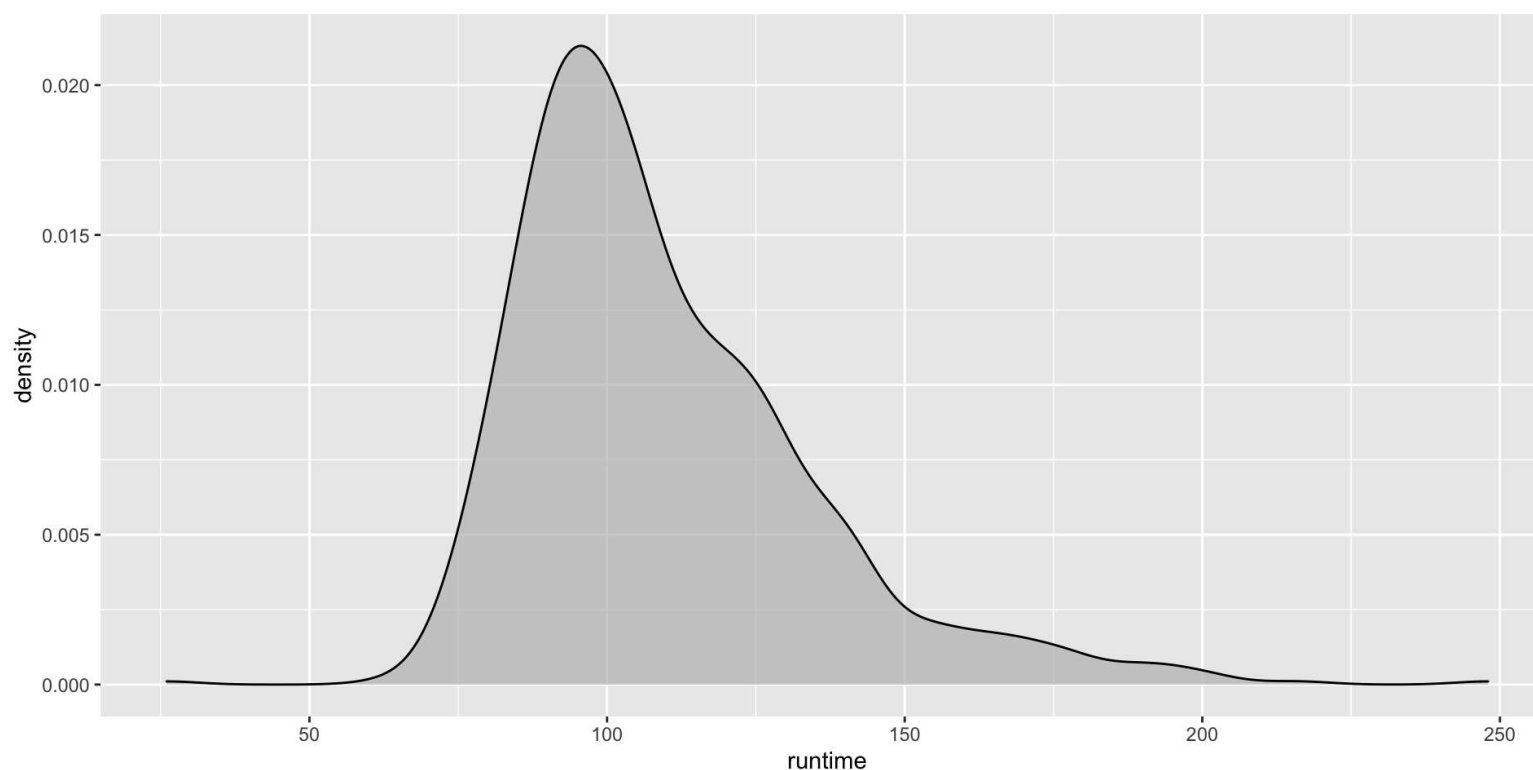https://github.com/ubco-mds-2022/Data-550

5

# Density plot

## Unlike Altair, ggplot has it's own density mark, …

```
1  ggplot(movies, aes(x = runtime)) +
2     geom_density(fill = 'grey', alpha = 0.7)
```



https://github.com/ubco-mds-2022/Data-550

6

# Unnesting the data

We need to unnest/explode on genres and countires.

```
1  free_genres <- movies %>% unnest(genres)
2  free_countries <- movies  %>%  unnest(countries)
3  free_both <- movies %>% unnest(genres) %>%  unnest(countries)
```

```
1  free_genres %>%
2    filter(, title ==  "All Dogs Go
3    select(genres, countries)
```

```
1  free_both %>%
2    filter(title ==  "All Dogs Go to
3    select(genres, countries)
```

| genres | countries |
| --- | --- |
| <chr> | <named list> |
| Fantasy | <chr [2]> |
| Animation | <chr [2]> |
| 2 rows | |

| genres | countries |
| --- | --- |
| <chr> | <chr> |
| Fantasy | United Kingdom |
| Fantasy | United States of America |
| Animation | United Kingdom |
| Animation | United States of America |
| 4 rows | |

https://github.com/ubco-mds-2022/Data-550

7

# Layered Density Plot
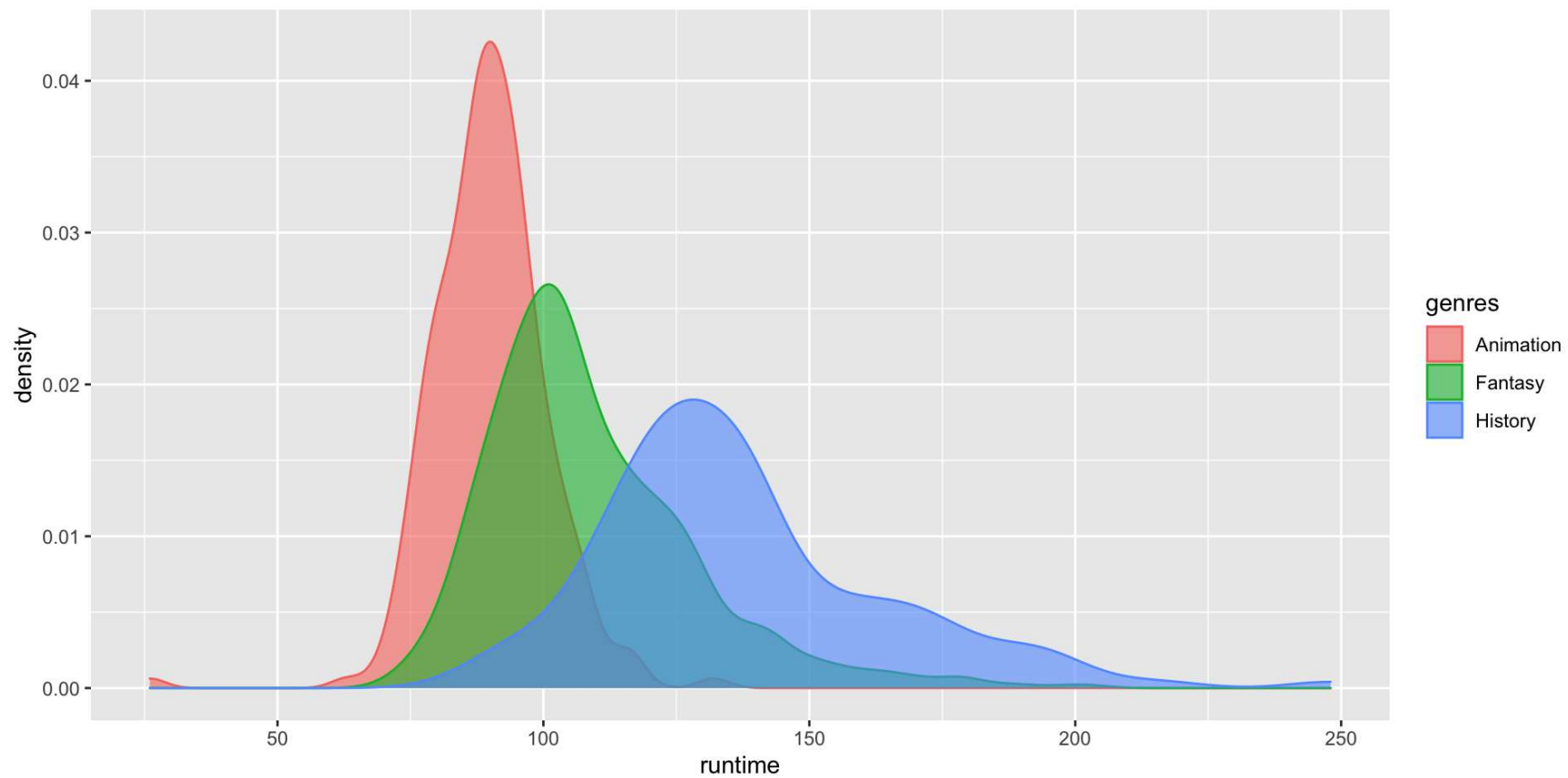
```
1  ggplot(free_genres, aes(x = runtime,
2         fill = genres,
3         color = genres)) +
4      geom_density(alpha = 0.6)
```

Notice how you can add the aesthetic rather than including it as an argument within `ggplot()`:

```
1  ggplot(free_genres) +
2      aes(x = runtime,
3          fill = genres,
4          color = genres) +
5      geom_density(alpha = 0.6)
```
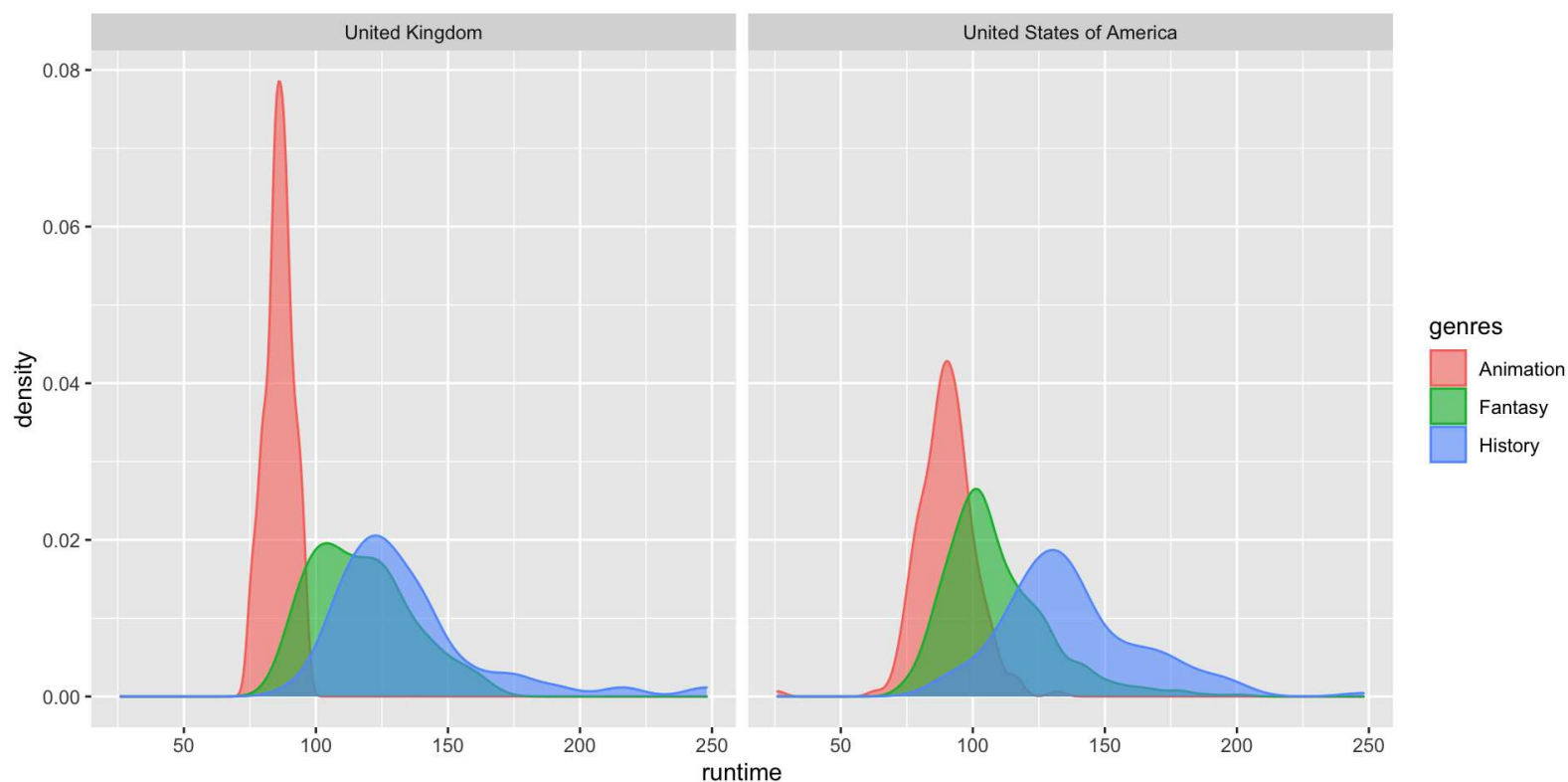
https://github.com/ubco-mds-2022/Data-550

8

# Layered Density Plot



https://github.com/ubco-mds-2022/Data-550

9

# Faceting
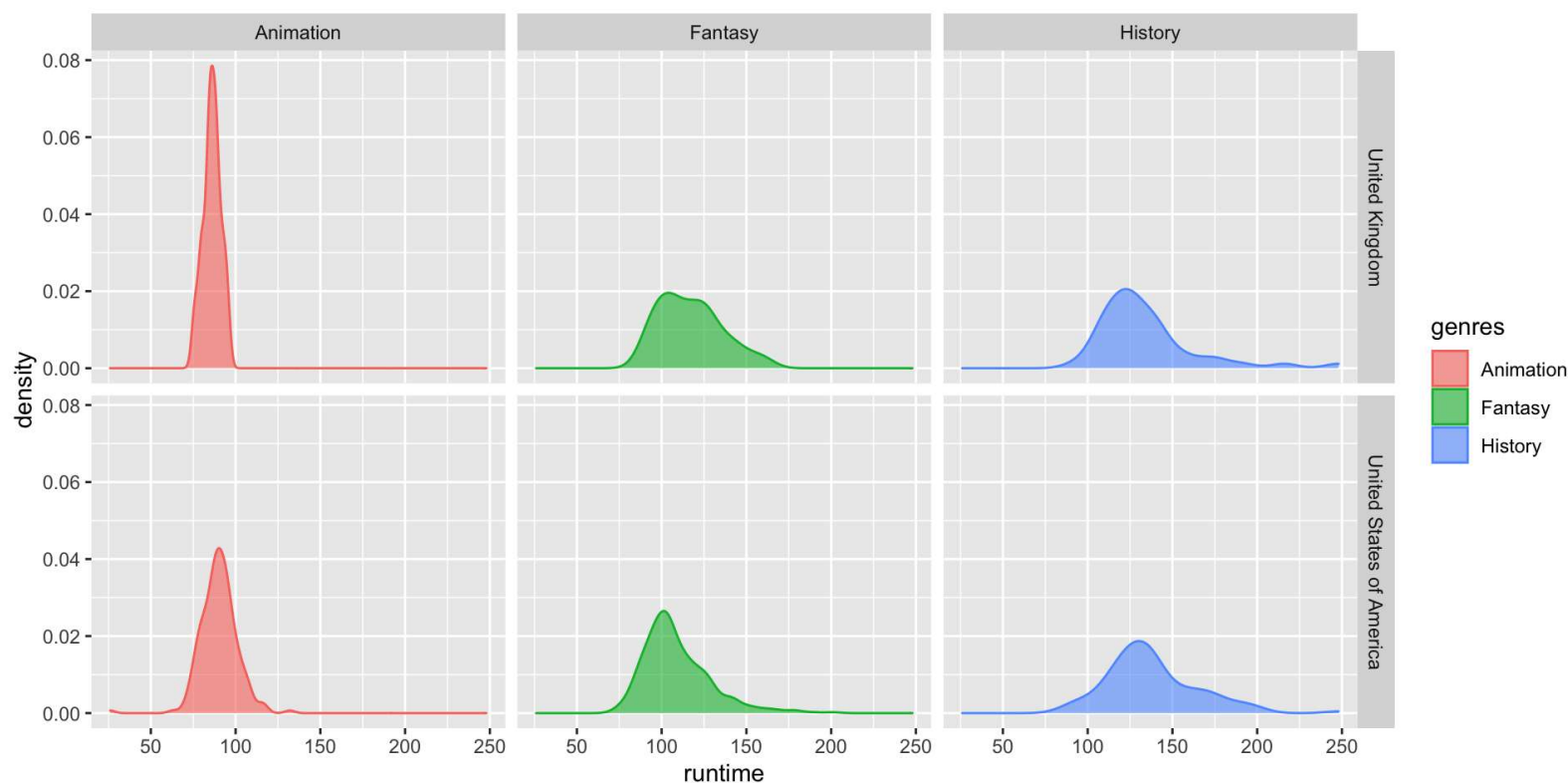
```
1  ggplot(free_both) +
2      aes(x = runtime, fill = genres, color = genres) +
3      geom_density(alpha = 0.6) +
4      facet_wrap(~countries)
```



https://github.com/ubco-mds-2022/Data-550

# Faceting (row and column)

```
1  ggplot(free_both, show.legend = FALSE) +
2      aes(x = runtime, fill = genres, color = genres) +
3      geom_density(alpha = 0.6) +
4      facet_grid(countries~genres)
```
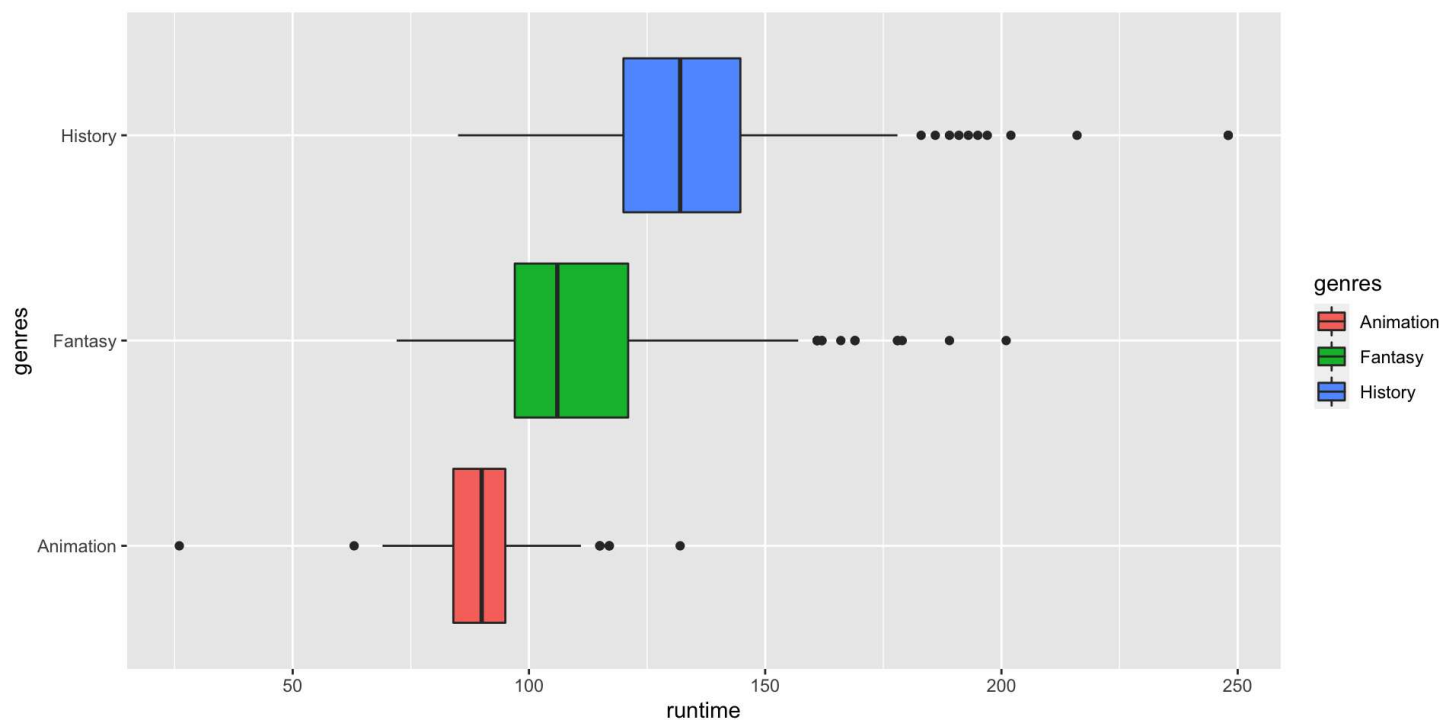


https://github.com/ubco-mds-2022/Data-550

11

# Boxplots

As in Altair, ggplot unsuprisingly has a boxplot geom, eg.

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = genres) +
3      geom_boxplot()
```
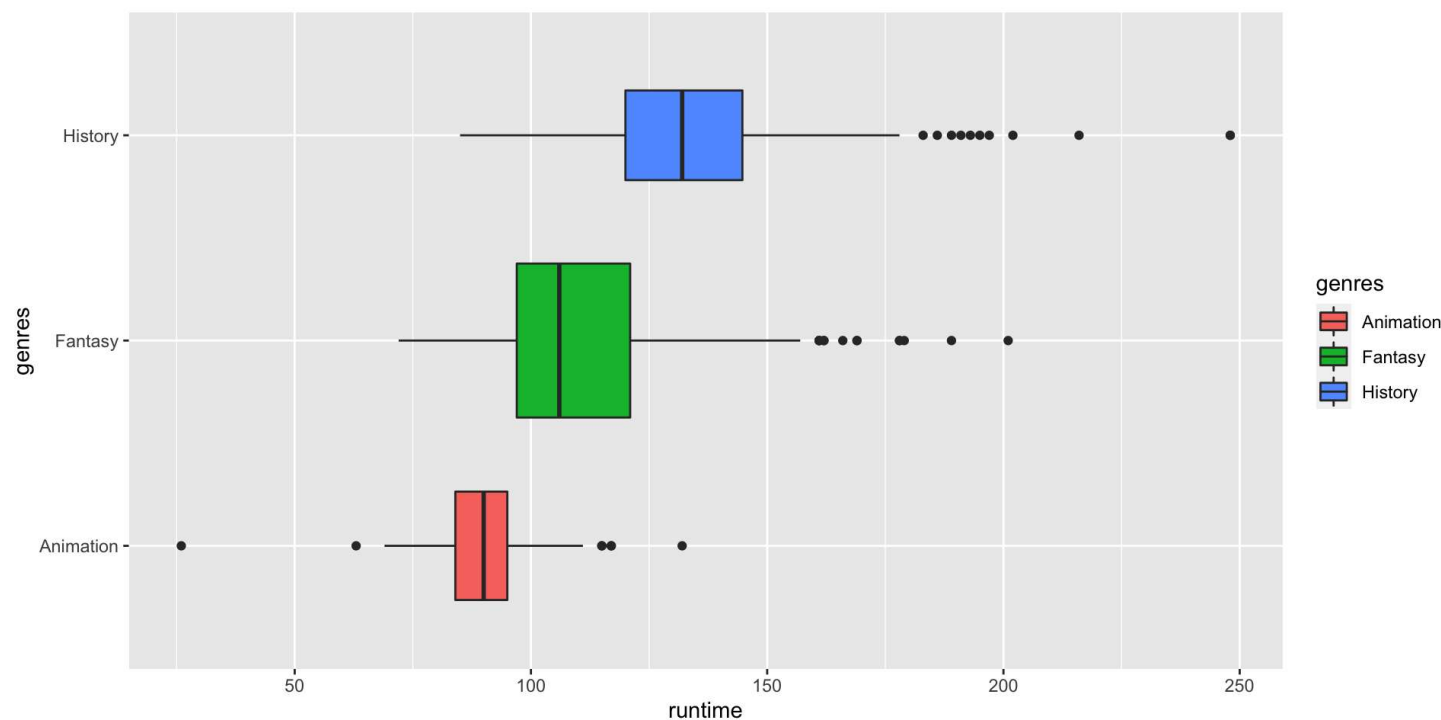


https://github.com/ubco-mds-2022/Data-550

12

# Scaled Boxplots

As in Altair, ggplot unsuprisingly has a boxplot geom, eg.

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres,fill = genres) +
3      geom_boxplot(varwidth = TRUE)
```
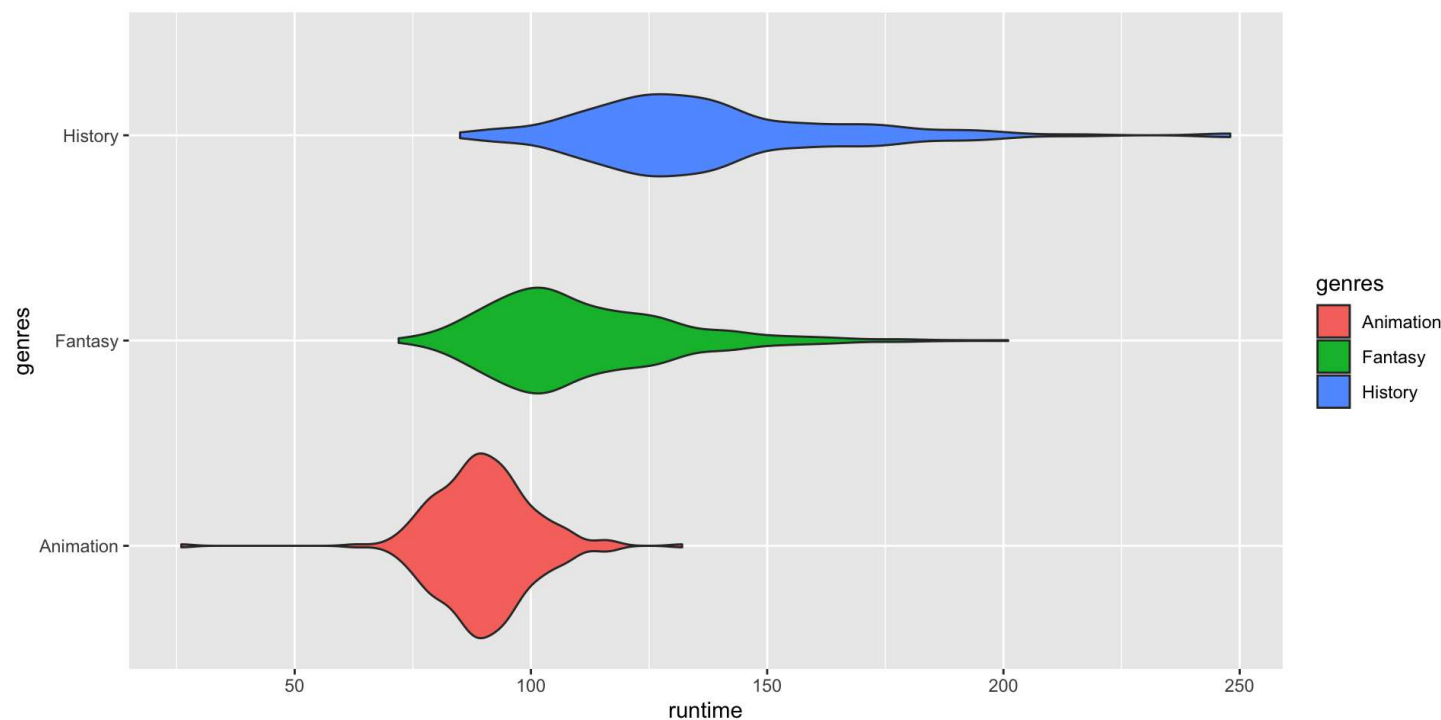


https://github.com/ubco-mds-2022/Data-550

13

# Violin Plots

## The change from boxplot to violin is extremely simple

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres,fill = genres) +
3      geom_violin()
```



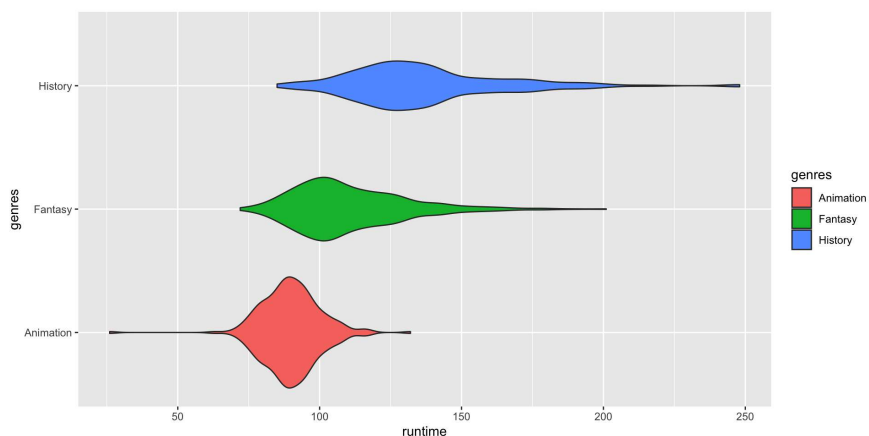https://github.com/ubco-mds-2022/Data-550

14

# What are violin plots

- Violin plots are similar to box plots, except that they also show the kernel probability density of the data at different values.

- Typically, violin plots will include a marker for the median of the data and a box indicating the interquartile range, as in standard box plots.

- The function `geom_violin()` is used to produce a violin plot.

https://github.com/ubco-mds-2022/Data-550
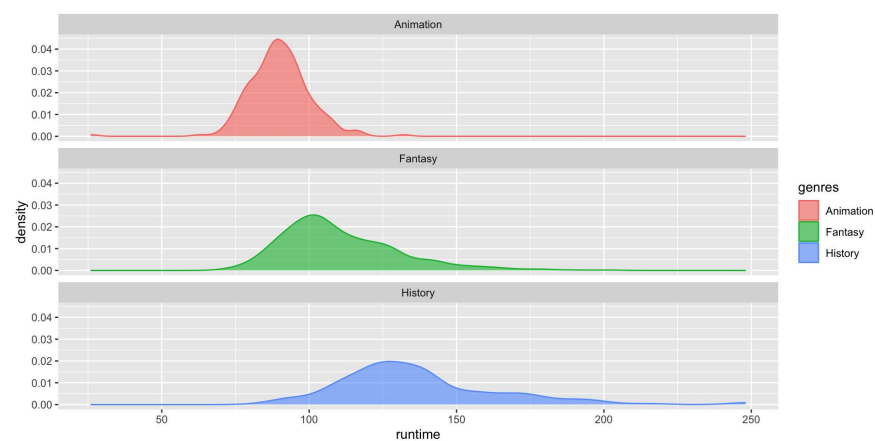
15

# Violin vs Faceted Density Plots

```
1  ggplot(free_both) +
2      aes(x = runtime,
3          y = genres,
4          fill = genres) +
5      geom_violin()
```

```
1  ggplot(free_both) +
2      aes(x = runtime, fill = genres
3      geom_density(alpha = 0.6) +
4      facet_wrap(~genres, ncol = 1)
```
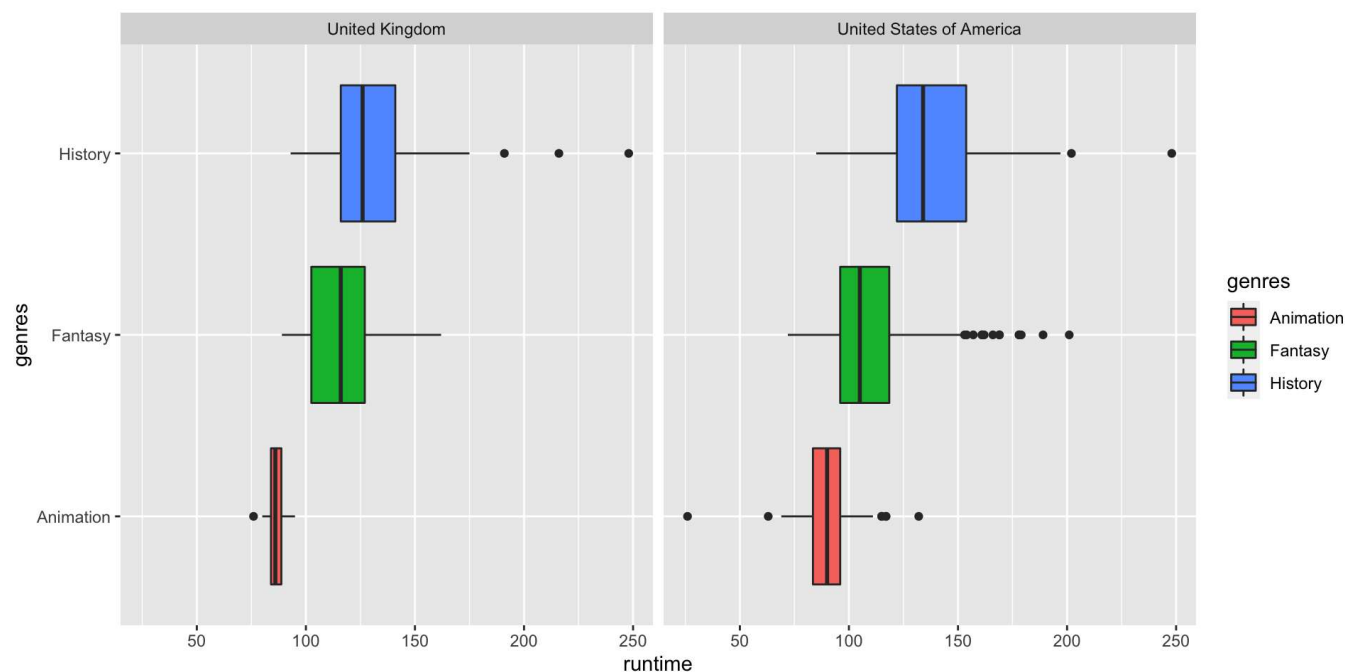




https://github.com/ubco-mds-2022/Data-550

16

# Faceted Boxplots

As with out density plots, we can also facet by country, eg.

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres,fill = genres) +
3      geom_boxplot() +
4      facet_wrap(~countries)
```
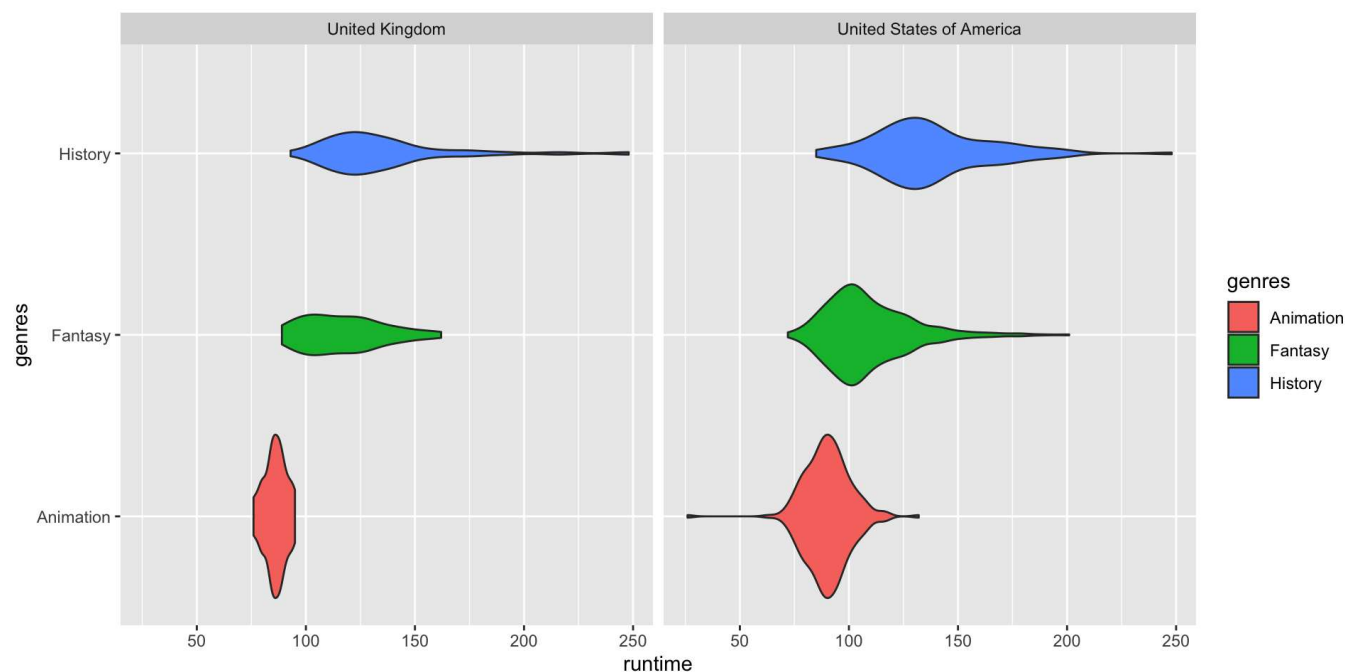


https://github.com/ubco-mds-2022/Data-550

17

# Violin Plots

## To get the violin plots, we simply change the geom:

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres,fill = genres) +
3      geom_violin() +
4      facet_wrap(~countries)
```
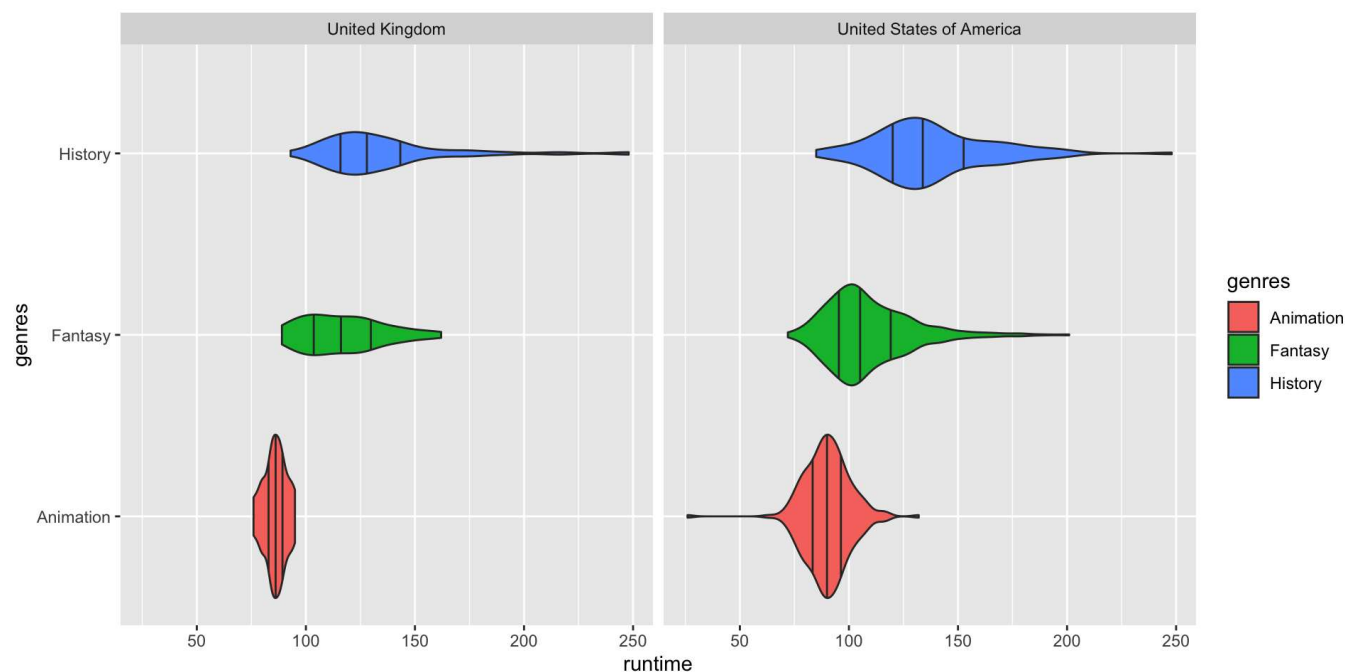


https://github.com/ubco-mds-2022/Data-550

18

# Layering Quanties

## We can layer the quantiles shown in the box plots

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = genres) +
3      geom_violin(draw_quantiles = c(0.25, 0.5, 0.75)) +
4      facet_wrap(~countries)
```



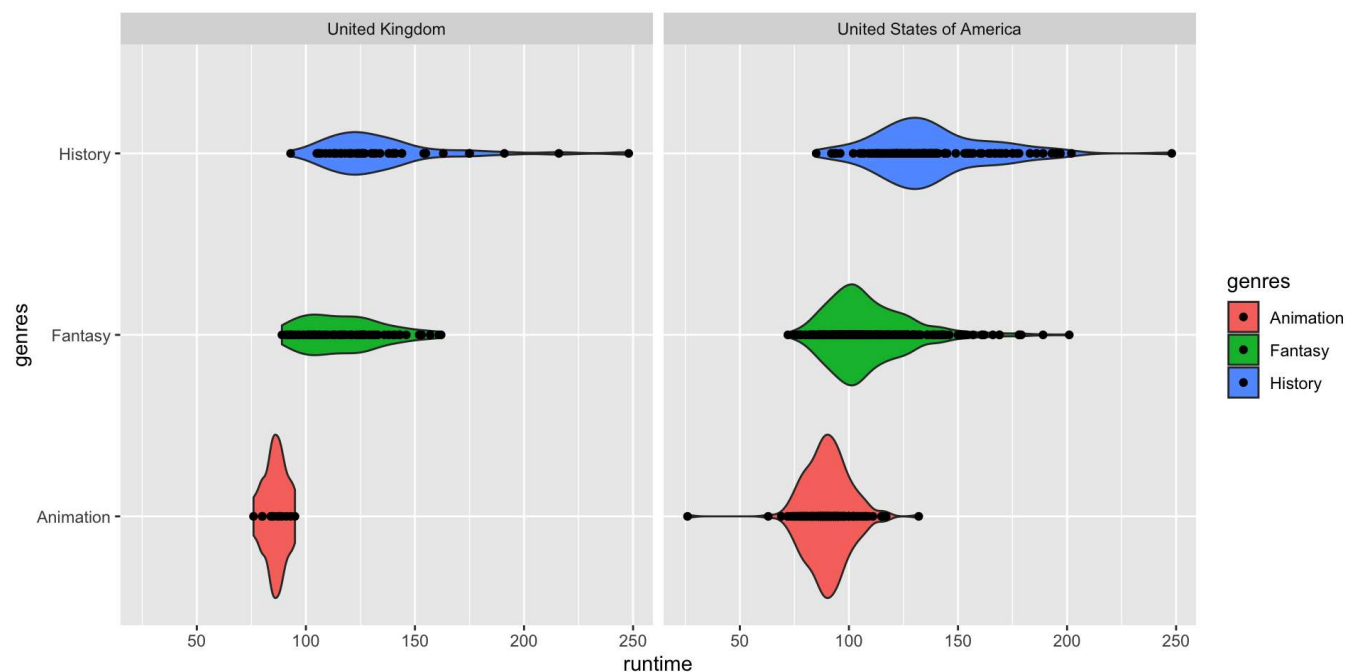https://github.com/ubco-mds-2022/Data-550

19

# Comments

- When possible, it is a good idea to have a look at where the individual data points are.

- Of course we could always layer on different marking of our data (using `geom_point()` for example)

- However when we have a lot of data, this could be impossible to read.

- For this we can use a categorical scatter plot where the dots are spread/jittered[1] randomly on the non-value axis so that they don't all overlap via `geom_jitter()`.

https://github.com/ubco-mds-2022/Data-550

20

1. jittering is not yet available in Altair

# Layering Points

## We can layer the points onto the violin plots:

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = genres) +
3      geom_violin() + geom_point() +
4      facet_wrap(~countries)
```
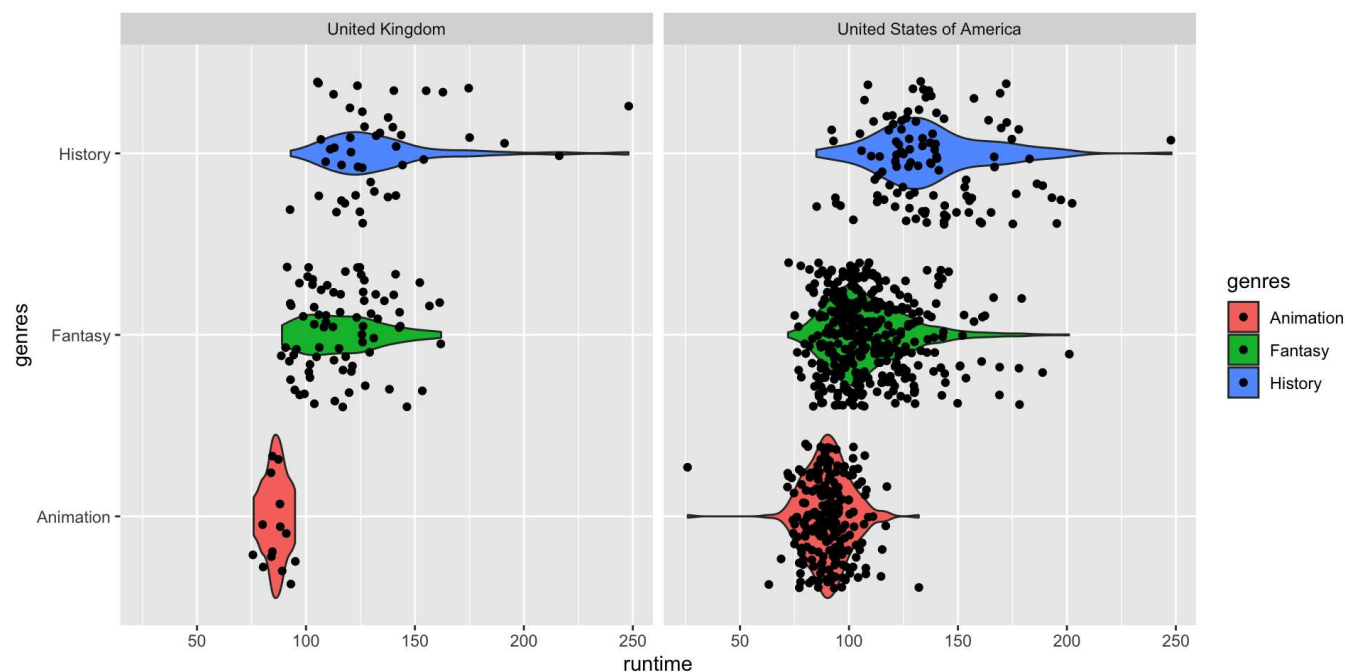


https://github.com/ubco-mds-2022/Data-550

21

# Jittering Data

"jittering" adds some noise to the location of each point

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = genres) +
3      geom_violin() + geom_jitter() +
4      facet_wrap(~countries)
```



https://github.com/ubco-mds-2022/Data-550
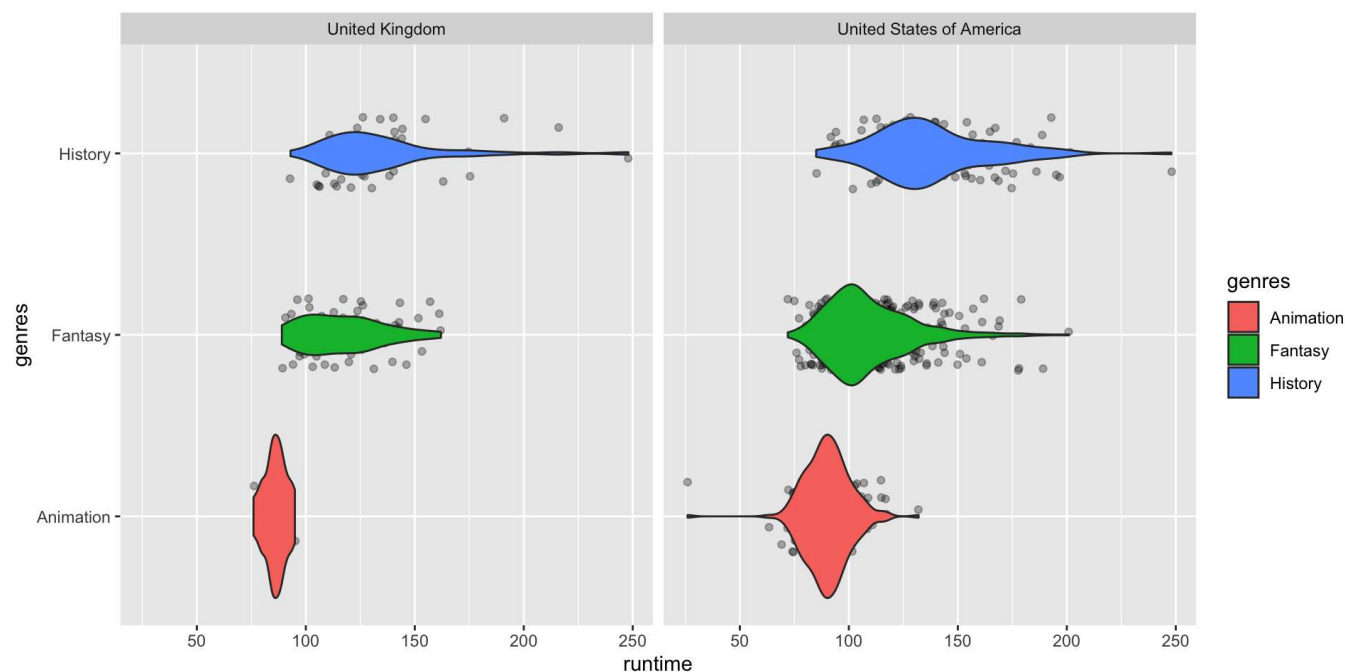
22

# Order matters

We can change the default height and order or layers

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = genres) +
3      geom_jitter(height = 0.2, alpha = 0.3) + geom_violin() +
4      facet_wrap(~countries)
```
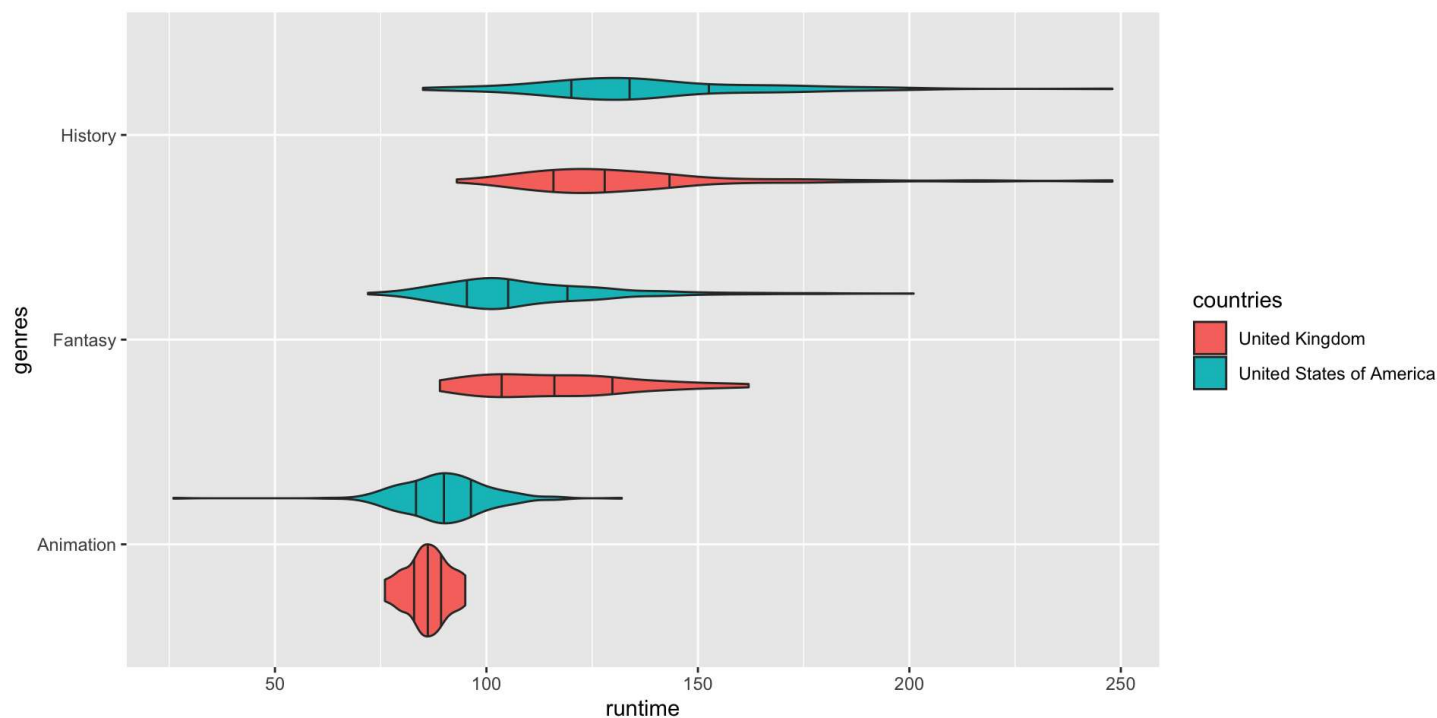


https://github.com/ubco-mds-2022/Data-550

23

# Unfaceting

Rather than faceting we could fill by countries

```
1  ggplot(free_both) +
2      aes(x = runtime, y = genres, fill = countries) +
3      geom_violin(draw_quantiles = c(0.25, 0.5, 0.75))
```



https://github.com/ubco-mds-2022/Data-550

24

https://github.com/ubco-mds-2022/Data-550