

# Tidy Tuesday Board Games

M. Zhang

2/21/2022

## Contents

<b>Introduction</b>	<b>1</b>
<b>Data Details</b>	<b>1</b>
Data Loading . . . . .	1
Data Dictionary . . . . .	2
<b>Data Exploration and wrangling</b>	<b>2</b>
Examining the raw data . . . . .	3
Joining two tables . . . . .	4
Checking for unique entries . . . . .	4
Grouping candidates . . . . .	6
Data cleaning for NA values . . . . .	7
<b>Data Visualizations</b>	<b>8</b>
<b>Conclusion</b>	<b>10</b>

## Introduction

We will be working with two datasets, **ratings** and **details**, sourced from BoardGameGeek via the Tidy Tuesday. This analysis will not be styled as a “final report”, but rather as a quick walk through of some data analysis and wrangling that took place while exploring the data.

The MSDSO Discord group for the University of Texas Masters in Data Science Online program will be doing weekly explorations of TidyTuesday as an exercise for improving their data science skill sets in a collaborative environment.

## Data Details

This dataset is pulled from the Tidy Tuesday Repository:

Thomas Mock (2022). Tidy Tuesday: A weekly data project aimed at the R ecosystem. <https://github.com/rfordatascience/tidytuesday>.

The original data is from Kaggle via Board Game Geek. The two data sets are joinable in the `id` column, and contain board game rating information, and board game details information.

## Data Loading

```
ratings <- readr::read_csv('https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/data/2022/2022-ratings.csv')
details <- readr::read_csv('https://raw.githubusercontent.com/rfordatascience/tidytuesday/master/data/2022/2022-details.csv')
```

## Data Dictionary

The following tables contain information in regards to the columns available.

### ratings.csv

variable	class	description
num	double	Game number
id	double	Game ID
name	character	Game name
year	double	Game year
rank	double	Game rank
average	double	Average rating
bayes_average	double	Bayes average rating
users Rated	double	Users rated
url	character	Game url
thumbnail	character	Game thumbnail

### details.csv

variable	class	description
num	double	Game number
id	double	Game ID
primary	character	Primary name
description	character	Description of game
yearpublished	double	Year published
minplayers	double	Min n of players
maxplayers	double	Max n of players
playingtime	double	Playing time in minutes
minplaytime	double	Min play time
maxplaytime	double	Max plat tome
minage	double	minimum age
boardgamecategory	character	Category
boardgamemechanic	character	Mechanic
boardgamefamily	character	Board game family
boardgameexpansion	character	Expansion
boardgameimplementation	character	Implementation
boardgamedesigner	character	Designer
boardgameartist	character	Artist
boardgamepublisher	character	Publisher
owned	double	Num owned
trading	double	Num trading
wanting	double	Num wanting
wishing	double	Num wishing

## Data Exploration and wrangling

Based on the data dictionary and information provided on the dataset, this appears to be a good opportunity to practice using `join` functions. Additionally, we will explore which columns are good candidates for

exploring summary data using the `group_by` function. In the end, we will attempt to put together a few meaningful graphics on the data summarizing user ratings and game rankings.

## Examining the raw data

First, we will use the `names` function to look at the information to quickly confirm the descriptions provided by the data dictionary.

```
ratings
```

```
# A tibble: 21,831 x 10
  num   id name      year rank average bayes_average users_rated url
  <dbl> <dbl> <chr>    <dbl> <dbl> <dbl>         <dbl>    <dbl> <chr>
1   105 30549 Pandemic    2008  106  7.59          7.49    108975 /boa~
2   189   822 Carcassonne 2000  190  7.42          7.31    108738 /boa~
3   428   13 Catan      1995  429  7.14          6.97    108024 /boa~
4    72 68448 7 Wonders    2010  73   7.74          7.63     89982 /boa~
5   103 36218 Dominion    2008  104  7.61          7.50     81561 /boa~
6   191  9209 Ticket to R~ 2004  192  7.41          7.30     76171 /boa~
7   100 178900 Codenames    2015  101  7.6           7.51     74419 /boa~
8    3 167791 Terraformin~ 2016   4   8.42          8.27     74216 /boa~
9   15 173346 7 Wonders D~ 2015  16   8.11          7.98     69472 /boa~
10  35  31260 Agricola    2007  36   7.93          7.81     66093 /boa~
# ... with 21,821 more rows, and 1 more variable: thumbnail <chr>
```

```
names(ratings)
```

```
[1] "num"          "id"           "name"         "year"
[5] "rank"         "average"      "bayes_average" "users_rated"
[9] "url"          "thumbnail"
```

```
ratings_count <- ratings %>% summarize(count = n())
```

```
details
```

```
# A tibble: 21,631 x 23
  num   id primary      description yearpublished minplayers maxplayers
  <dbl> <dbl> <chr>    <chr>         <dbl>         <dbl>    <dbl>
1    0 30549 Pandemic    In Pandemi~    2008           2          4
2    1   822 Carcassonne Carcassonn~    2000           2          5
3    2   13 Catan      In CATAN (~    1995           3          4
4    3 68448 7 Wonders    You are th~    2010           2          7
5    4 36218 Dominion    &quot;You ~    2008           2          4
6    5  9209 Ticket to Ride With elega~    2004           2          5
7    6 178900 Codenames    Codenames ~    2015           2          8
8    7 167791 Terraforming Ma~ In the 240~    2016           1          5
9    8 173346 7 Wonders Duel In many wa~    2015           2          2
10   9  31260 Agricola    Descriptio~    2007           1          5
# ... with 21,621 more rows, and 16 more variables: playingtime <dbl>,
# minplaytime <dbl>, maxplaytime <dbl>, minage <dbl>,
# boardgamecategory <chr>, boardgamemechanic <chr>, boardgamefamily <chr>,
# boardgameexpansion <chr>, boardgameimplementation <chr>,
# boardgamedesigner <chr>, boardgameartist <chr>, boardgamepublisher <chr>,
# owned <dbl>, trading <dbl>, wanting <dbl>, wishing <dbl>
```

```
names(details)
```

```
[1] "num"          "id"
```

```

[3] "primary"          "description"
[5] "yearpublished"    "minplayers"
[7] "maxplayers"       "playingtime"
[9] "minplaytime"      "maxplaytime"
[11] "minage"           "boardgamecategory"
[13] "boardgamemechanic" "boardgamefamily"
[15] "boardgameexpansion" "boardgameimplementation"
[17] "boardgamedesigner" "boardgameartist"
[19] "boardgamepublisher" "owned"
[21] "trading"          "wanting"
[23] "wishing"

```

```
details_count <- details %>% summarize(count = n())
```

- Ratings table total entries: 21831
- Details table total entries: 21631

## Joining two tables

Based off the first 6 entries shown in each data table, it does appear that the `id` column should allow for linkage of the two tables. There are less entries in the details dataset though, so we will left join on the details dataset and drop non-matching entries.

```
board_games <- left_join(details, ratings, by = "id")
board_games
```

```

# A tibble: 21,631 x 32
  num.x   id primary      description yearpublished minplayers maxplayers
  <dbl> <dbl> <chr>          <chr>          <dbl>         <dbl>      <dbl>
1     0 30549 Pandemic    In Pandemi~    2008           2          4
2     1  822 Carcassonne Carcassonn~    2000           2          5
3     2   13 Catan      In CATAN (~    1995           3          4
4     3 68448 7 Wonders  You are th~    2010           2          7
5     4 36218 Dominion  &quot;You ~    2008           2          4
6     5  9209 Ticket to Ride With elega~    2004           2          5
7     6 178900 Codenames  Codenames ~    2015           2          8
8     7 167791 Terraforming Ma~ In the 240~    2016           1          5
9     8 173346 7 Wonders Duel  In many wa~    2015           2          2
10    9 31260 Agricola    Descriptio~    2007           1          5
# ... with 21,621 more rows, and 25 more variables: playingtime <dbl>,
#   minplaytime <dbl>, maxplaytime <dbl>, minage <dbl>,
#   boardgamecategory <chr>, boardgamemechanic <chr>, boardgamefamily <chr>,
#   boardgameexpansion <chr>, boardgameimplementation <chr>,
#   boardgamedesigner <chr>, boardgameartist <chr>, boardgamepublisher <chr>,
#   owned <dbl>, trading <dbl>, wanting <dbl>, wishing <dbl>, num.y <dbl>,
#   name <chr>, year <dbl>, rank <dbl>, average <dbl>, bayes_average <dbl>, ...

```

We can see based off the total number of entries in the new table `board_games` that we were successful in our `left_join`

## Checking for unique entries

Next, we will see if each game (`primary`) has a unique entry in the data set.

```

board_games %>%
  group_by(primary) %>%
  summarize(

```

```

n = n()) %>%
  arrange(desc(n))

# A tibble: 21,236 x 2
  primary          n
  <chr>          <int>
1 Robin Hood      6
2 Chaos           4
3 Cosmic Encounter 4
4 Gangster        4
5 Gettysburg      4
6 Saga            4
7 Warhammer 40,000: Kill Team 4
8 Airlines        3
9 Ali Baba        3
10 Around the World in 80 Days 3
# ... with 21,226 more rows

```

It can be seen that there are multiple entries per game. Based off the data set information, it may be that the same game with different publication dates maintain separate entries. We will modify the original summary to check:

```

board_games %>%
  group_by(primary, yearpublished) %>%
  summarize(
    n = n(), .groups = 'drop') %>%
  arrange(desc(n))

# A tibble: 21,626 x 3
  primary          yearpublished    n
  <chr>          <dbl> <int>
1 "Cahoots"      2018     2
2 "Chaos"        2010     2
3 "Columbus"     1991     2
4 "DIG"          2017     2
5 "Loch Ness"    2010     2
6 "'65: Squad-Level Combat in the Jungles of Vietnam" 2016     1
7 "'CA' Tactical Naval Warfare in the Pacific, 1941-45" 1973     1
8 "'Wacht am Rhein': The Battle of the Bulge, 16 Dec 44 - ~ 1977     1
9 "\"La Garde recule!\"" 2011     1
10 "\"Oh My God! There\\'s An Axe In My Head.\" The Game of~ 2014     1
# ... with 21,616 more rows

```

We still see some duplicate entries so we will pull up a specific game to examine the data:

```

board_games %>%
  select(primary, boardgamepublisher, average, description) %>%
  filter(primary == "Cahoots")

# A tibble: 2 x 4
  primary boardgamepublisher average description
  <chr>    <chr>                <dbl> <chr>
1 Cahoots ['Gamewright', 'Oya', 'Brain Picnic', 'Lifestyle ~ 7.18 Cooperatio~
2 Cahoots ['Mayday Games'] 6.47 Welcome to~

```

It appears that the games have separate and therefore have separate entries. As a crude solution for the purposes of this analysis, will move forward since the number of entries in which multiples appear are small.

## Grouping candidates

For summarizing and visualizing the data, a few potential grouping candidates based off the `details` dataset are `primary`, `boardgamecategory`, `boardgamemechanic`, and `boardgamepublisher`. We will use `sapply` and `n_distinct`.

```
board_games %>%
  select(primary, boardgamecategory, boardgamemechanic, boardgamepublisher) %>%
  sapply(n_distinct)
```

primary	boardgamecategory	boardgamemechanic	boardgamepublisher
21236	6731	8292	11266

Unfortunately, we can see from the above that the number of unique entries for each row is much higher than anticipated. To understand why, we will look at the actual data contained in these columns.

```
board_games %>%
  select(
    primary,
    boardgamecategory,
    boardgamemechanic,
    boardgamepublisher) %>%
  head
```

```
# A tibble: 6 x 4
  primary boardgamecategory boardgamemechan~ boardgamepublis~
  <chr>    <chr>                  <chr>          <chr>
1 Pandemic ['Medical']          ['Action Points~ ['Z-Man Games',~
2 Carcassonne ['City Building', 'Medieval'~ ['Area Majority~ ['Hans im Glück~
3 Catan      ['Economic', 'Negotiation'] ['Dice Rolling'~ ['KOSMOS', '999~
4 7 Wonders  ['Ancient', 'Card Game', 'Ci~ ['Drafting', 'H~ ['Repos Product~
5 Dominion   ['Card Game', 'Medieval']   ['Deck, Bag, an~ ['Rio Grande Ga~
6 Ticket to Ride ['Trains']              ['Card Drafting~ ['Days of Wonde~
```

It appears that the `boardgame` columns all actually contain more than single value entries, leading to large number of permutations of values. These columns will not be suitable for quick data analysis. Thus, we will examine some other parameters.

```
board_games %>%
  select(playingtime, minplayers, maxplayers, minage) %>%
  sapply(n_distinct)
```

playingtime	minplayers	maxplayers	minage
119	11	52	21

```
board_games %>%
  select(
    playingtime,
    minplayers,
    maxplayers,
    minage) %>%
  head
```

```
# A tibble: 6 x 4
  playingtime minplayers maxplayers minage
  <dbl>      <dbl>      <dbl> <dbl>
1      45         2         4      8
2      45         2         5      7
3     120         3         4     10
```

4	30	2	7	10
5	30	2	4	13
6	60	2	5	8

```
board_games %>%
  select(
    playingtime,
    minplayers,
    maxplayers,
    minage) %>%
  summary()
```

playingtime	minplayers	maxplayers	minage
Min. : 0.00	Min. : 0.000	Min. : 0.00	Min. : 0.000
1st Qu.: 25.00	1st Qu.: 2.000	1st Qu.: 4.00	1st Qu.: 8.000
Median : 45.00	Median : 2.000	Median : 4.00	Median :10.000
Mean : 90.51	Mean : 2.007	Mean : 5.71	Mean : 9.611
3rd Qu.: 90.00	3rd Qu.: 2.000	3rd Qu.: 6.00	3rd Qu.:12.000
Max. :60000.00	Max. :10.000	Max. :999.00	Max. :25.000

These are all numerical values with reasonable distributions (though there appears to be some outlier values which may be removed later).

## Data cleaning for NA values

```
board_games %>%
  summarize(across(everything(), ~ sum(is.na(.)))) %>%
  tidyr::pivot_longer(everything()) %>%
  arrange(desc(value)) %>%
  deframe()
```

boardgameimplementation	boardgameexpansion	boardgameartist
16769	16125	5907
boardgamefamily	boardgamemechanic	boardgamedesigner
3761	1590	596
boardgamecategory	thumbnail	description
283	6	1
boardgamepublisher	num.x	id
1	0	0
primary	yearpublished	minplayers
0	0	0
maxplayers	playingtime	minplaytime
0	0	0
maxplaytime	minage	owned
0	0	0
trading	wanting	wishing
0	0	0
num.y	name	year
0	0	0
rank	average	bayes_average
0	0	0
usersRated	url	
0	0	

NA counting methodology was taken from stackexchange

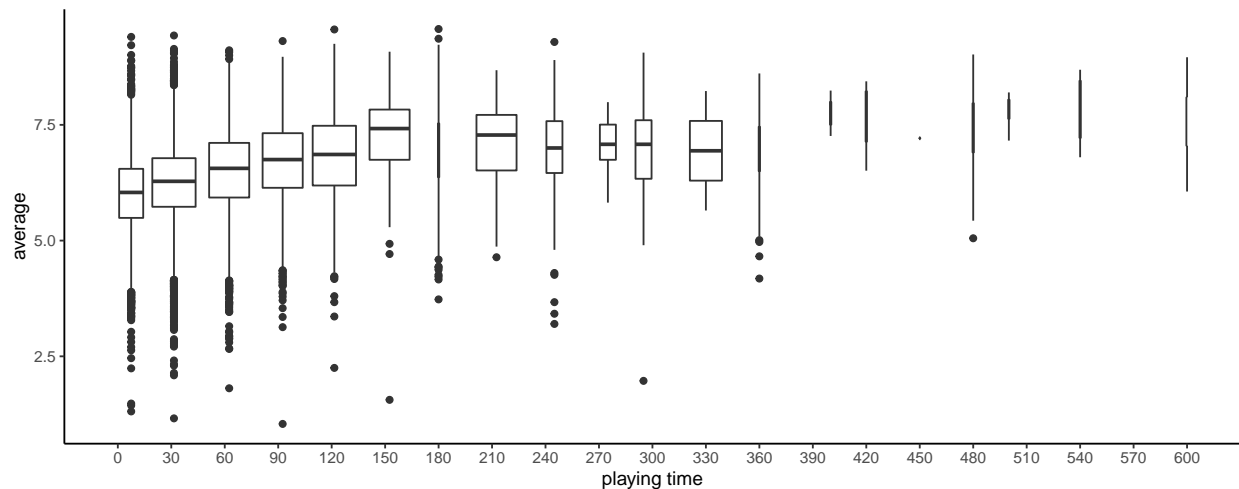
We see there are a number of NA values, but that `playingtime`, `rank`, and `average` are all populated. Thus, we will now move onto data visualizations.

## Data Visualizations

```
board_games %>%
  select(
    average,
    rank,
    playingtime,
    minplayers,
    maxplayers,
    minage) %>%
  ggplot() +
  aes(x = playingtime, y = average, group = cut_width(playingtime, 30)) +
  geom_boxplot() +
  scale_x_continuous(
    name = "playing time",
    breaks = seq(0, 600, 30),
    labels = seq(0, 600, 30),
    limits = c(0, 600)) +
  theme_classic()
```

Warning: Removed 192 rows containing missing values (stat\_boxplot).

Warning: Removed 1 rows containing missing values (geom\_segment).



```
board_games %>%
  select(
    average,
    rank,
    playingtime,
    minplayers,
    maxplayers,
    minage) %>%
  ggplot() +
  aes(playingtime) +
  geom_histogram()
```



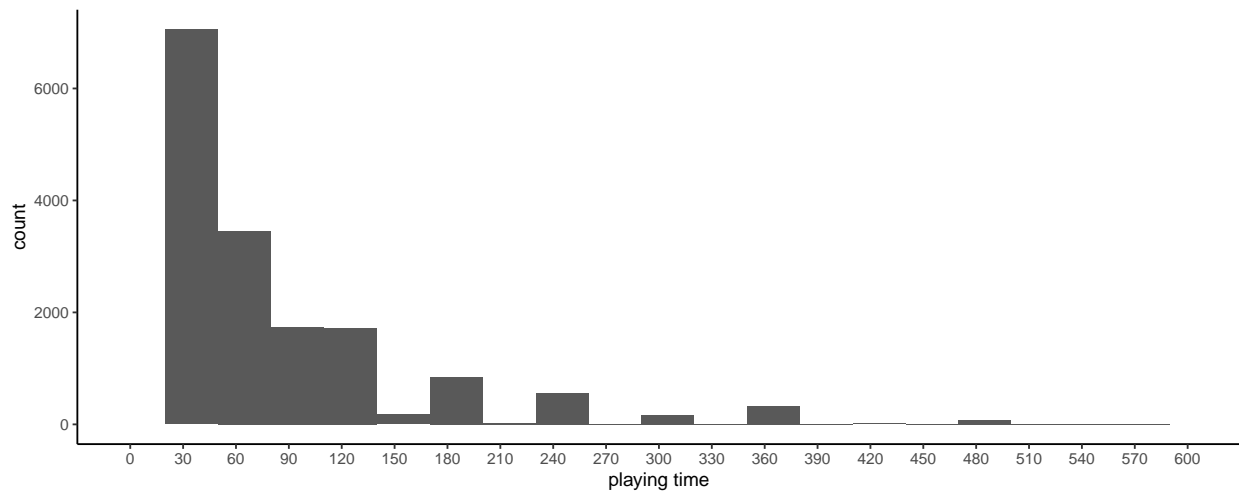
```

    binwidth = 30,
    center = 5) +
scale_x_continuous(
  name = "playing time",
  breaks = seq(0, 600, 30),
  labels = seq(0, 600, 30),
  limits = c(0, 600)) +
theme_classic()

```

Warning: Removed 192 rows containing non-finite values (stat\_bin).

Warning: Removed 2 rows containing missing values (geom\_bar).

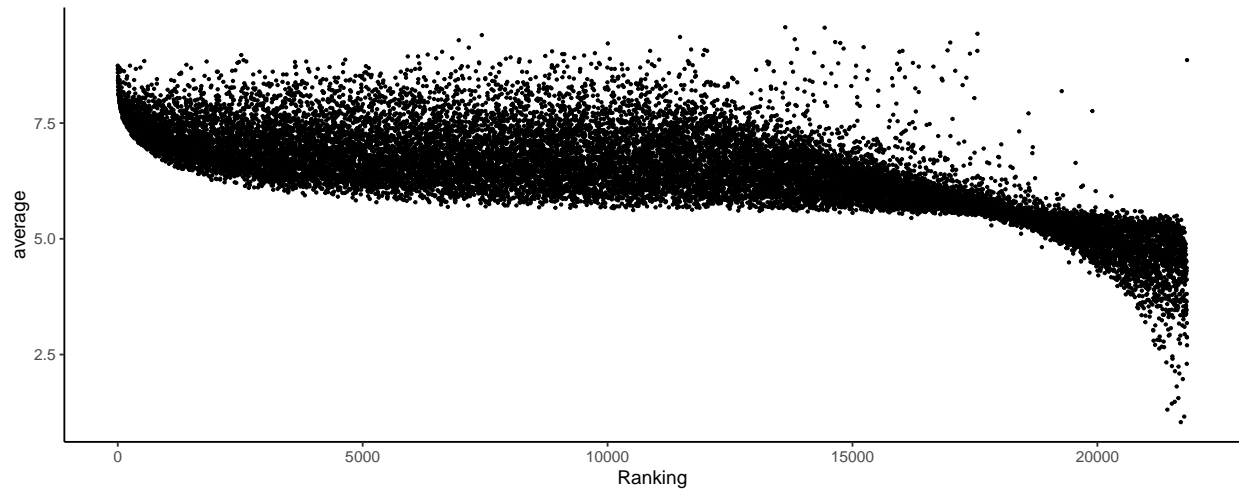


The range of scores appears to diminish as the playing time for the game increases. It can also be seen that the number of games which extend to longer playtimes tapers off very quickly. There were also more rows of data being dropped than expected - this is a point for future analysis if this dataset is revisited.

```

board_games %>%
  select(
    average,
    rank,
    playingtime,
    minplayers,
    maxplayers,
    minage) %>%
  ggplot() +
  aes(x = rank, y = average) +
  geom_point(size = 0.5) +
  scale_x_continuous(
    name = "Ranking") +
  theme_classic()

```



Ranking and score have relationship where the second order rate changes signs at larger numerical values for ranking.

## Conclusion

This ended up being a shorter initial exploration of the board games dataset. One finding of note was that there was a lot of game information which still needs to be parsed out of the `boardgame` columns which actually contain delimited lists of information which could allow for more interesting analysis and charts in the future.

Thanks for reading and joining me on my journey to improve my skills in data analysis.