

# **Capstone Project 2: Milestone Report**

## **Board Game Review**

### **Problem Statement**

Board games have regained popularity in recent years, they make it possible to entertain many afternoons with family or evenings with friends during hours. There are more and more of them and different kinds of games ranging from collaborative games to strategy games and more family games.

**The purpose of this project is to construct a model using machine learning and natural language processing to predict the rating of board games based on the reviews of clients and some other features as the number of players, the average time of a game, the number of rates...**

Potential clients of this project could be the board game designers so that they know the favorite game characteristics of the clients or they could be the managers of gaming shops, board game bar or larger department stores so that they know which games they will be most sensitive to sell.

For that, I will use scraping and APIs on the website: <https://boardgamegeek.com> following commands on this website: [https://boardgamegeek.com/wiki/page/BGG\\_XML\\_API2](https://boardgamegeek.com/wiki/page/BGG_XML_API2). The website [boardgamegeek.com](https://boardgamegeek.com) is a reference in terms of board games to receive advices and explanations from a large community.

### **Description of the dataset**

We extracted from the fifty most rated boardgames on boardgamegeek.com, the following features for each of these boardgames:

- The ID
- The name
- The year of design
- The minimum number of players required
- The maximum number of players required
- The minimum number of minutes required to complete the game
- The maximum number of minutes required to complete the game
- The minimum age required to play
- The category

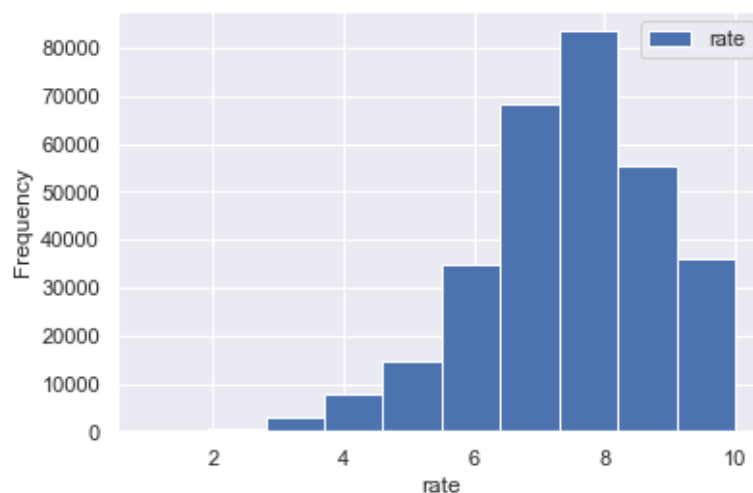
- The number of rates
- All the reviews about the game with the username of the person who wrote it and the associated note score out of 10.

Then, we registered the new dataframe in csv format.

In a new notebook, we loaded the dataframe, we shuffled it because it is sorting by games and then, we deleted the rows containing missing values in the 'review' column. The final dataset contains 304864 reviews.

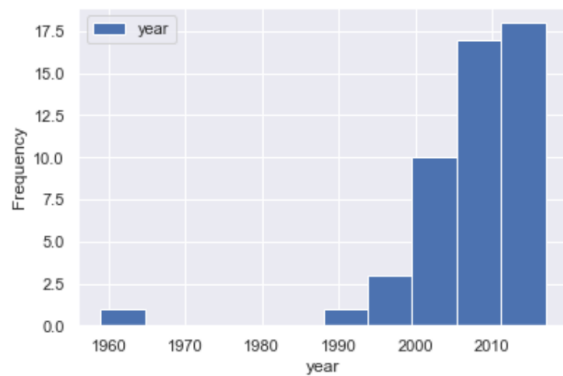
## Visual exploratory data analysis

The following histogram shows us the distribution of all the rates gave by players for the 50 more rated games of the website boardgamegeek.com. We can see a normal distribution of the rates distributed around 7.



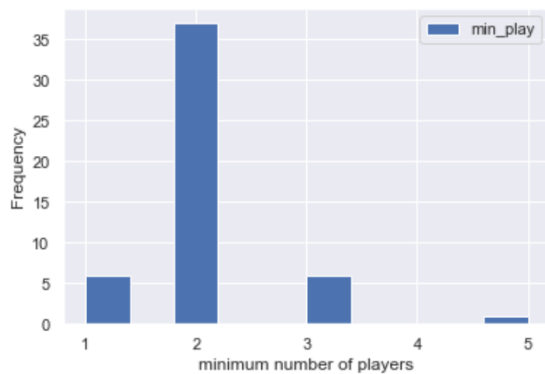
Then, I created a new dataframe with one boardgame per row and the average of all the rates given for each game in the 'rate' column using 'groupby' method. From this new dataframe, we have drawn several plots.

## Distribution of years of design of the boardgames



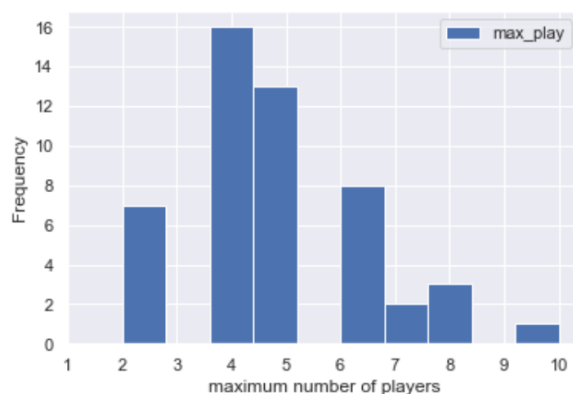
This histogram shows us that most of the games have been created during the last 20 years. There is one outlier game created in 1960.

## Distribution of the minimum number of players required for a boardgame



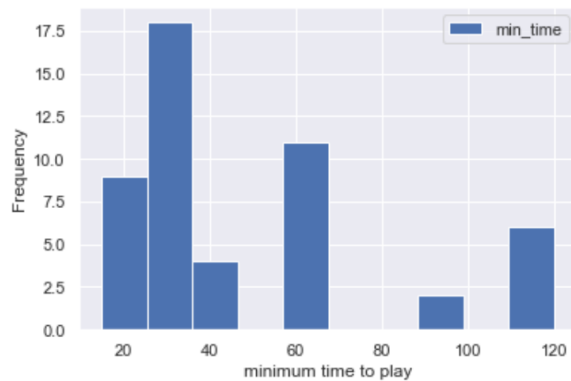
We can see that the distribution of the minimum number of players is between 1 and 5 with a big majority of 2 players.

## Distribution of the maximum number of players required for a boardgame



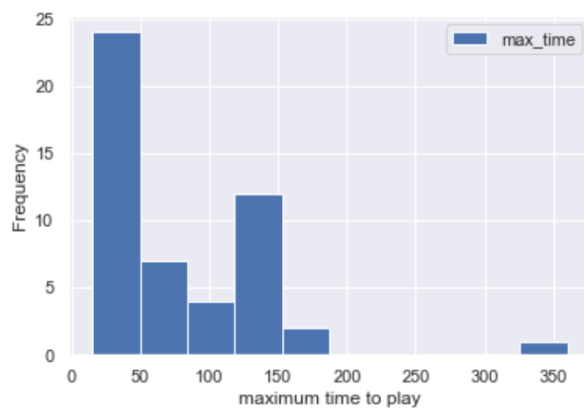
We can see that the distribution of the maximum number of players is between 2 and 10 with a majority of 4 players.

### Distribution of the minimum number of minutes required to complete a game



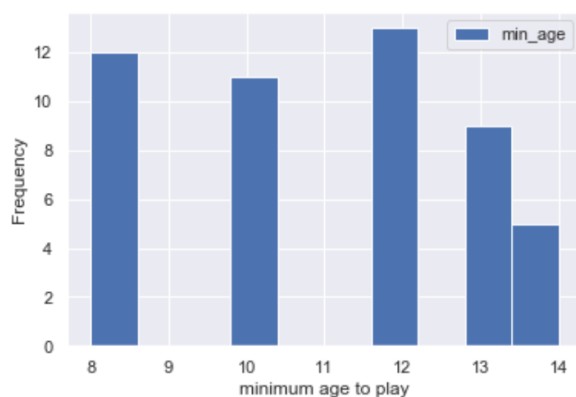
This plot shows us that the distribution of the minimum number of minutes for a play is between 20 and 120 minutes with a majority of 30 minutes.

### Distribution of the maximum number of minutes required to complete a game



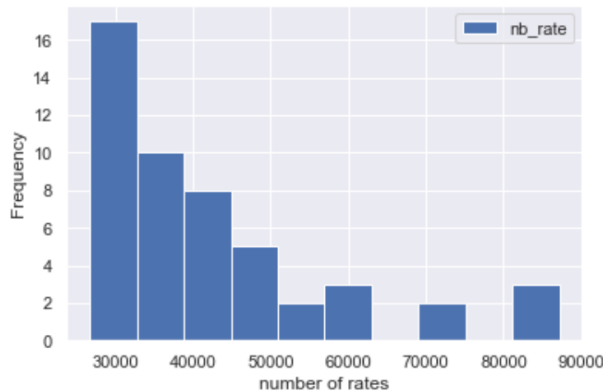
This plot shows us that the distribution of the maximum number of minutes for a play is most often less than 50 minutes.

### Distribution of the minimum of age required to play



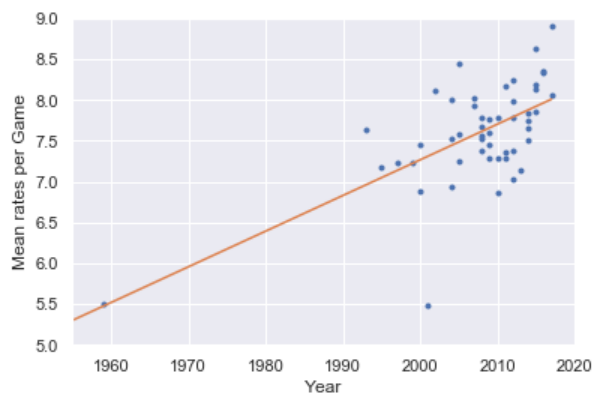
This histogram shows us that the distribution of the minimum age required to play is between 8 and 14 years old.

## Distribution of the number of rates by boardgame



This histogram shows us that the distribution of the number of rates by game is between 30000 and 85000 rates.

## Relationship between the mean rates per boardgames and the years of the design

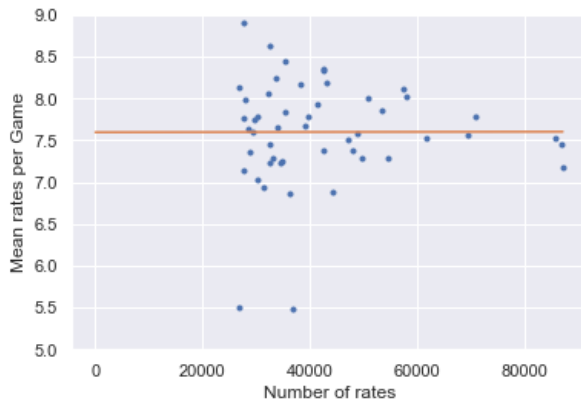


We can see that seems to have a positive correlation between the year of the design of a game and the mean of rates of this game. More a boardgame is recent and more he seems to have higher mean rates.



We can observe on the previous plot an outlier value which is a very old boardgame compared to the others (designed before 1960). We removed this value in this second plot to compare the slop without the outlier value. We can see that the slop is similar to the previous plot, it seems to have a positive correlation between the year of the design of a game and his average rate.

## Relationship between the mean rates per boardgame and the number of rates per games



We can see that seems to have no correlation between the number of the rates of a game and the mean of rates of this game. So, even if we chose to take the 50 most rated games from boardgamegeek.com for this project, it doesn't seem to be a bias regarding the rating of these 50 games.

## Statistical exploratory data analysis

```
=====
                        OLS Regression Results
=====
Dep. Variable:          rate      R-squared:                0.636
Model:                  OLS      Adj. R-squared:             0.575
Method:                 Least Squares      F-statistic:          10.48
Date:                   Sun, 22 Sep 2019    Prob (F-statistic):    1.56e-07
Time:                   11:25:49           Log-Likelihood:       -21.989
No. Observations:       50              AIC:                  59.98
Df Residuals:           42              BIC:                  75.27
Df Model:               7
Covariance Type:        nonrobust
=====
                        coef      std err          t      P>|t|      [0.025      0.975]
-----
Intercept      -79.2237       14.211      -5.575      0.000     -107.904     -50.544
year              0.0428        0.007        6.055      0.000         0.029         0.057
min_play       -0.0935        0.102     -0.914      0.366     -0.300         0.113
max_play       -0.0989        0.038     -2.587      0.013     -0.176     -0.022
min_time        0.0043        0.003        1.347      0.185     -0.002         0.011
max_time       -0.0003        0.002     -0.167      0.868     -0.004         0.003
min_age         0.0947        0.032        2.920      0.006         0.029         0.160
nb_rate         8.415e-06    3.98e-06        2.116      0.040      3.9e-07    1.64e-05
=====
Omnibus:            21.055    Durbin-Watson:           2.033
Prob(Omnibus):      0.000    Jarque-Bera (JB):        39.160
Skew:               -1.223    Prob(JB):                3.14e-09
Kurtosis:           6.579    Cond. No.                1.11e+07
=====
```

We can see that the features `min_play`, `min_time` and `max_time` have both a very high `p_value`. They are not good predictors for the mean rate of boardgames. The other features have a significant `p_value` (for an  $\alpha=0.05$ ). The `year`, the `min_age` and the `nb_rate` have positive coefficients, the mean rate of boardgames increases when the value of these three features also increase. The `max_play` has negative coefficient, the mean rates of boardgames increases when the value of this feature decreases. The `year` and the `min_age` have the lower `p_values`, less than

0.01. If the year increases by one unit (1 year), the mean rates of boardgames increases by 0.04.  
If the min\_age increases by one unit (1 year), the mean rates of boardgames increases by 0.095.