

# Google Play Store Data Analysis Report

Suheng Yao

2024-11-18

## Abstract

The main purpose of this report is about prediction of rating scores of the apps in google play store and find out the important variables that can influence the rating score. The method used are multiple linear regression model and linear mixed effect model. After testing for the model results, the multiple linear regression gets the better performance, and the important variables are number of reviews, price, number of installs, size, category and the time the app was last updated. In the discussion section, the report had a literature review, comparing the approach used in this report with methods used by other researchers. Also, some biases within the dataset used is mentioned in the discussion section.

## Introduction

Google Play Store has always been the most popular and largest app store for Android phone users across the world. Since it is pre-installed on supporting Android devices, and the operative system holds over 70 percent of the global market, Google Play Store becomes the default app hub for most of the Android users worldwide. Up until the second quarter of 2024, there are 2.26 million apps available in the store[1].

In this project, the main question of interest is: What are the important factors that can affect the apps ratings? The reasons why I am interested in this question are divided into two part: first of all, since I use Google Play to install apps every day, I am always curious about what kind of apps can have good ratings scores and why some apps I find easy to use get low rating score. Secondly, by doing this analysis, if I want to build apps to publish on Google Play in the future, then I will know which essential factors to focus on or the popular fields to go into.

This project report will be divided into four parts: methods section will talk about model selection and model building; results section will talk about the important findings after fitting the models; discussion section will further analyze those findings and talk about the next steps for the analysis; appendix is the last part, which will include EDA and initial understandings of the data.

## Data Cleaning

```
## Rows: 10,841
## Columns: 13
## $ App      <chr> "Photo Editor & Candy Camera & Grid & ScrapBook", "Colo~
## $ Category <chr> "ART_AND_DESIGN", "ART_AND_DESIGN", "ART_AND_DESIGN", "~
## $ Rating   <dbl> 4.1, 3.9, 4.7, 4.5, 4.3, 4.4, 3.8, 4.1, 4.4, 4.7, 4.4, ~
## $ Reviews  <chr> "159", "967", "87510", "215644", "967", "167", "178", "~
## $ Size     <chr> "19M", "14M", "8.7M", "25M", "2.8M", "5.6M", "19M", "29~
## $ Installs <chr> "10,000+", "500,000+", "5,000,000+", "50,000,000+", "10~
```

```
## $ Type      <chr> "Free", "Free", "Free", "Free", "Free", "Free", "Free", ~
## $ Price     <chr> "0", "0", "0", "0", "0", "0", "0", "0", "0", "0", "0", ~
## $ Content.Rating <chr> "Everyone", "Everyone", "Everyone", "Teen", "Everyone", ~
## $ Genres    <chr> "Art & Design", "Art & Design;Pretend Play", "Art & Des~
## $ Last.Updated <chr> "January 7, 2018", "January 15, 2018", "August 1, 2018"~
## $ Current.Ver <chr> "1.0.0", "2.0.0", "1.2.4", "Varies with device", "1.1", ~
## $ Android.Ver <chr> "4.0.3 and up", "4.0.3 and up", "4.0.3 and up", "4.2 an~
```

Looking at the structure of the dataset, there are in total 10,841 observations and 13 variables related to each app. The response variable is **Rating**, which is a continuous variable. By looking at the column names of the data, “Category” and “Genre” seem to have similar data, I can print out the unique values in those two columns:

```
## [1] "ART_AND_DESIGN"      "AUTO_AND_VEHICLES"    "BEAUTY"
## [4] "BOOKS_AND_REFERENCE" "BUSINESS"              "COMICS"
## [7] "COMMUNICATION"      "DATING"                "EDUCATION"
## [10] "ENTERTAINMENT"      "EVENTS"                "FINANCE"
## [13] "FOOD_AND_DRINK"     "HEALTH_AND_FITNESS"   "HOUSE_AND_HOME"
## [16] "LIBRARIES_AND_DEMO" "LIFESTYLE"             "GAME"
## [19] "FAMILY"             "MEDICAL"               "SOCIAL"
## [22] "SHOPPING"           "PHOTOGRAPHY"           "SPORTS"
## [25] "TRAVEL_AND_LOCAL"   "TOOLS"                  "PERSONALIZATION"
## [28] "PRODUCTIVITY"       "PARENTING"              "WEATHER"
## [ reached getOption("max.print") -- omitted 4 entries ]

## [1] "Art & Design"          "Art & Design;Pretend Play"
## [3] "Art & Design;Creativity" "Art & Design;Action & Adventure"
## [5] "Auto & Vehicles"      "Beauty"
## [7] "Books & Reference"    "Business"
## [9] "Comics"               "Comics;Creativity"
## [11] "Communication"        "Dating"
## [13] "Education;Education"  "Education"
## [15] "Education;Creativity" "Education;Music & Video"
## [17] "Education;Action & Adventure" "Education;Pretend Play"
## [19] "Education;Brain Games" "Entertainment"
## [21] "Entertainment;Music & Video" "Entertainment;Brain Games"
## [23] "Entertainment;Creativity" "Events"
## [25] "Finance"              "Food & Drink"
## [27] "Health & Fitness"     "House & Home"
## [29] "Libraries & Demo"     "Lifestyle"
## [ reached getOption("max.print") -- omitted 90 entries ]
```

Based on the printed results above, “Genre” is just the more detailed classification of “Category”, since in this project, I mostly focus on the general groups of apps, I can remove the “Genre” column and change the Category name to lower case.

```
##                               App      Category Rating Reviews
## 1 Photo Editor & Candy Camera & Grid & ScrapBook art_and_design    4.1    159
## 2                               Coloring book moana art_and_design    3.9    967
##   Size Installs Type Price Content.Rating      Last.Updated Current.Ver
## 1   19M   10,000+ Free     0      Everyone January 7, 2018      1.0.0
## 2   14M  500,000+ Free     0      Everyone January 15, 2018      2.0.0
##   Android.Ver
```

```
## 1 4.0.3 and up
## 2 4.0.3 and up
## [ reached 'max' / getOption("max.print") -- omitted 4 rows ]
```

Also, from the distinct values of category, there is a category called “1.9”, which does not make sense. Since there is only one record of data with category “1.9”, I can just remove this record from the dataset.

Now, let’s check if there are any duplicated apps in the dataset:

```
## [1] 1181
```

From the analysis above, there are 1181 duplicated apps in the dataset, to make further analysis easier, I will just keep the first occurrence of each app record.

Now, I need to check if there are NA values in the dataset:

```
##      App      Category      Rating      Reviews
## Length:9659 Length:9659 Min. :1.000 Length:9659
## Class :character Class :character 1st Qu.:4.000 Class :character
##      Size      Installs      Type      Price
## Length:9659 Length:9659 Length:9659 Length:9659
## Class :character Class :character Class :character Class :character
## Content.Rating Last.Updated Current.Ver Android.Ver
## Length:9659 Length:9659 Length:9659 Length:9659
## Class :character Class :character Class :character Class :character
## [ reached getOption("max.print") -- omitted 5 rows ]
```

From the result shown above, in the response variable “Rating”, there are 1463 missing values. To maintain the size of the data, I will fill in those missing values using median values of the “Rating” column.

```
##      App      Category      Rating      Reviews
## Length:9659 Length:9659 Min. :1.000 Length:9659
## Class :character Class :character 1st Qu.:4.000 Class :character
##      Size      Installs      Type      Price
## Length:9659 Length:9659 Length:9659 Length:9659
## Class :character Class :character Class :character Class :character
## Content.Rating Last.Updated Current.Ver Android.Ver
## Length:9659 Length:9659 Length:9659 Length:9659
## Class :character Class :character Class :character Class :character
## [ reached getOption("max.print") -- omitted 4 rows ]
```

Since it makes more sense to talk about installs, reviews, size and price in numerical values, I will change those variables from categorical to numeric:

```
##      App      Category      Rating      Reviews
## Length:9659 Length:9659 Min. :1.000 Min. : 0
## Class :character Class :character 1st Qu.:4.000 1st Qu.: 25
## Size(in MB) Installs      Type      Price
## Min. : 1.00 Length:9659 Length:9659 Min. : 0.000
## 1st Qu.: 5.10 Class :character Class :character 1st Qu.: 0.000
## Content.Rating Last.Updated Current.Ver Android.Ver
## Length:9659 Length:9659 Length:9659 Length:9659
## Class :character Class :character Class :character Class :character
## [ reached getOption("max.print") -- omitted 5 rows ]
```

Since Most of the NA values in Size are related to “Varies with Device” value in the original dataset, and in this project, I want to mostly focus on the app with fixed size, I will just remove those app records with varied app sizes.

Transforming Installs is a more complex matter, I will solve this problem differently. First check the distinct values in the variable “Installs”:

```
## [1] "10,000+"      "500,000+"      "5,000,000+"     "50,000,000+"
## [5] "100,000+"     "50,000+"       "1,000,000+"     "10,000,000+"
## [9] "5,000+"       "100,000,000+"  "1,000+"         "500,000,000+"
## [13] "50+"          "100+"          "500+"           "10+"
## [17] "1+"           "5+"            "1,000,000,000+" "0+"
```

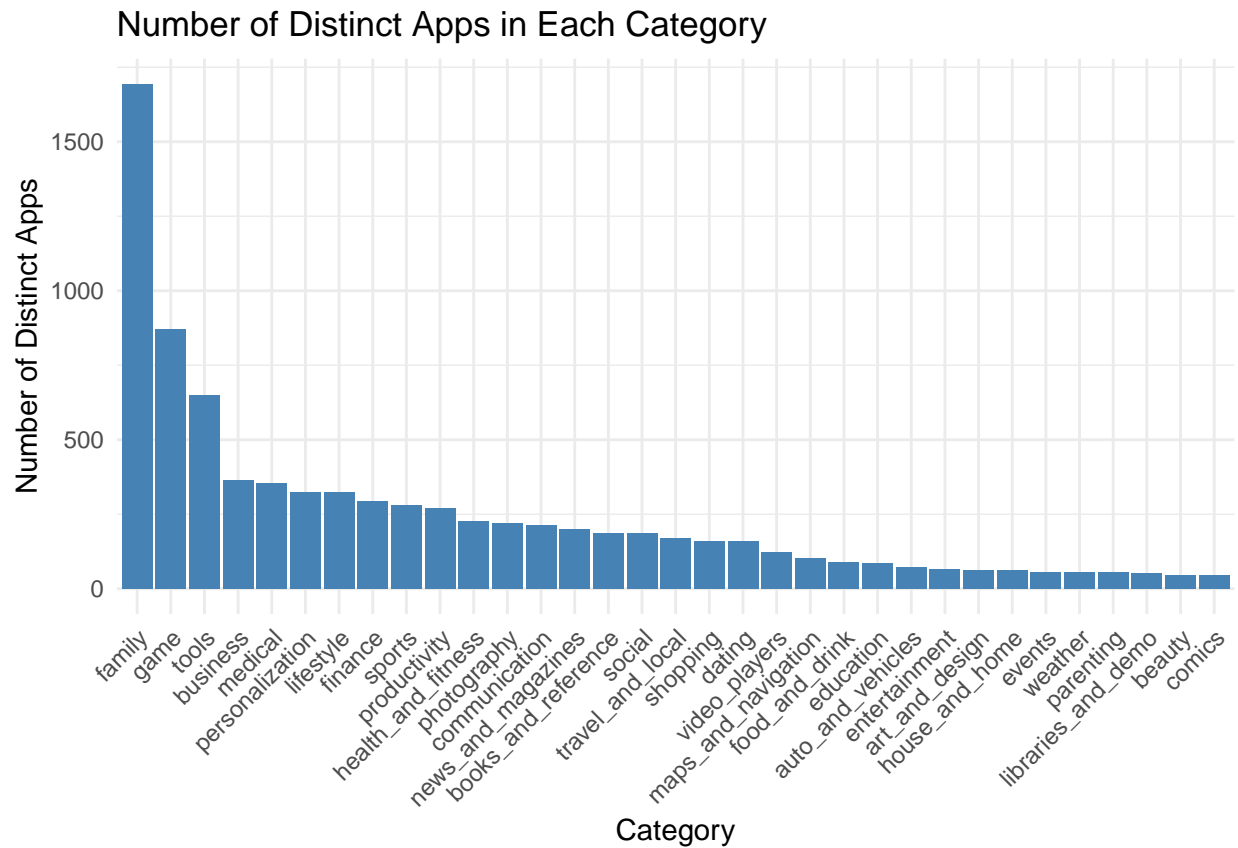
To convert those values into numeric values and avoid duplicate values, take “500+” as an example, I will change it to a value between 500 to 1000 because “1000+” will be on a different level.

```
##      App      Category      Rating      Reviews
## Length:8118 Length:8118 Min. :1.000 Min. : 0
## Class :character Class :character 1st Qu.:4.000 1st Qu.: 17
## Size(in MB) Installs      Type      Price
## Min. : 1.00 Min. :0.000e+00 Length:8118 Min. : 0.000
## 1st Qu.: 5.10 1st Qu.:2.821e+03 Class :character 1st Qu.: 0.000
## Content.Rating Last.Updated Current.Ver Android.Ver
## Length:8118 Length:8118 Length:8118 Length:8118
## Class :character Class :character Class :character Class :character
## [ reached getOption("max.print") -- omitted 4 rows ]
```

After those data cleaning is done, I will start doing EDA.

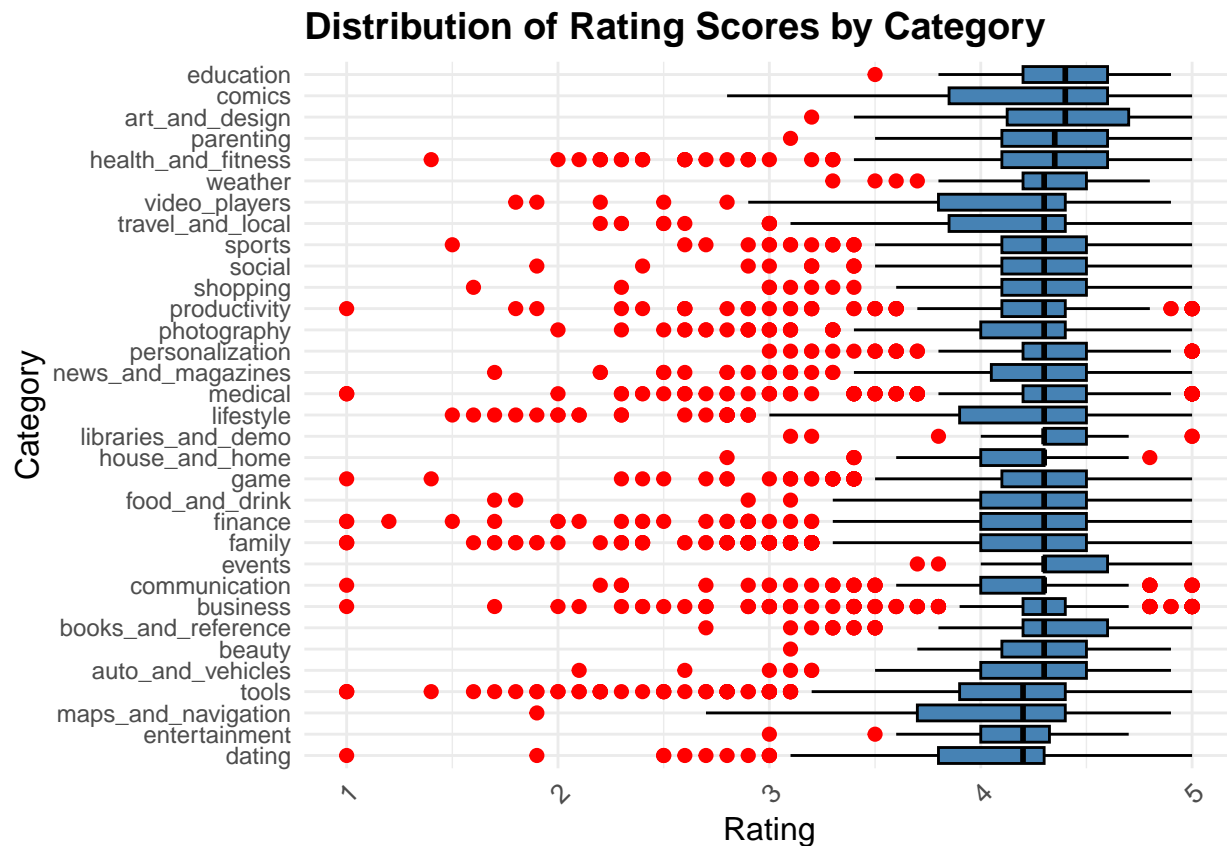
## EDA

### Step 1: Bar Plot Counting Distinct Apps in Each Category



From the count of distinct apps, it is clear that category family gets the most number of apps, with more than 1500 apps in the category. The second most category is game, which is only half the number of the family category, and the third most is the tools category. The category with the least number of distinct apps is comics, with only less than 250 apps.

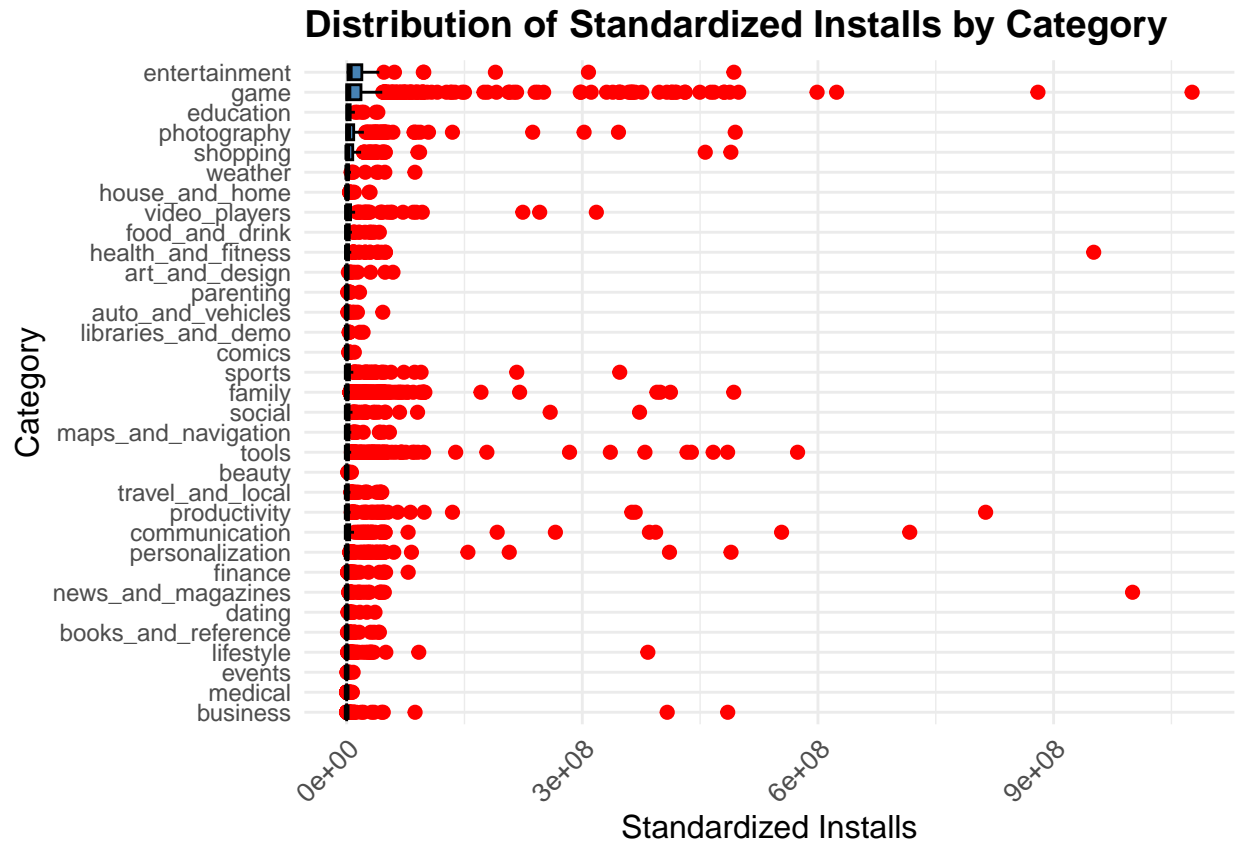
## Step 2: Box Plot of Distribution of Rating Scores for Each Category



From the box plot for each category above, it is clear that almost all the categories have median value greater than 4. Also, maybe due to less number of apps in comics category, it is the only category without outliers. What's more, the last four categories at bottom: Tools, Maps and Navigation, Entertainment and Dating tend to have lower median values compared to other groups, which may reflect the user dissatisfaction for those app groups.

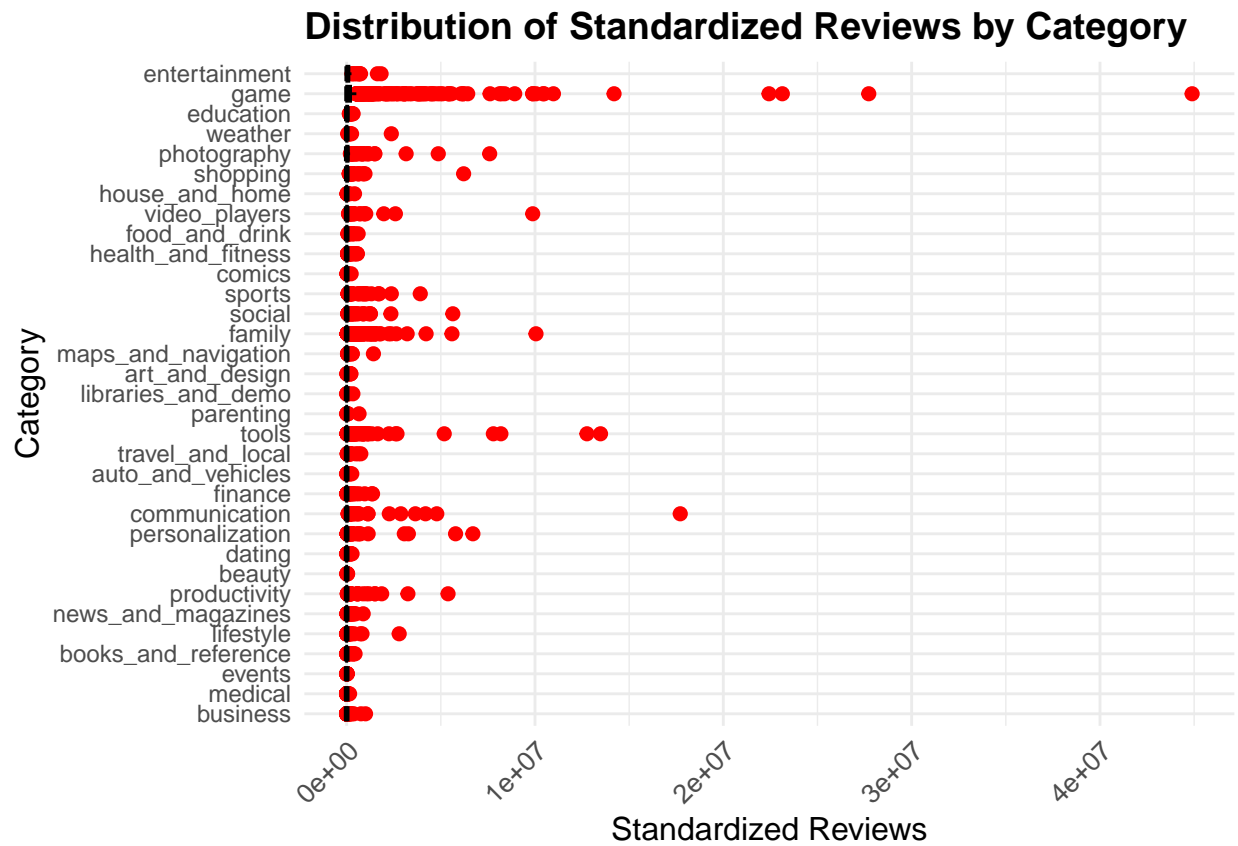
## Step 3: Table and Box Plot of Distribution of Number of Installs in Each Category

```
## # A tibble: 33 x 5
##   Category      count    min  median    max
##   <chr>      <int> <dbl>   <dbl>   <dbl>
## 1 family      1694     0  74861 492808674
## 2 game         870     1 2297828. 1076347463
## 3 tools        651     1  56676 574082114
## 4 business     365     0    971 485125907
## 5 medical      353     0   3492 7318962
## 6 personalization 325     0  34699 489348676
## 7 lifestyle     323     0  19204 383449407
## 8 finance       293     0  34503 78351602
## 9 sports        281     3  76462 347653545
## 10 productivity 271     0  44067 813489255
## # i 23 more rows
```



The table is based on original values of installs, and to make every numeric on a similar scale, I try to standardize the install variable, and the box plot is based on standardized install values. From this boxplot above, in general, games tend to have the more number of installs, and the range of app installs in the game category also tend to be greater than other categories. However, there is one outlier in personalization category, which has over  $1.5 * 10^9$  installs.

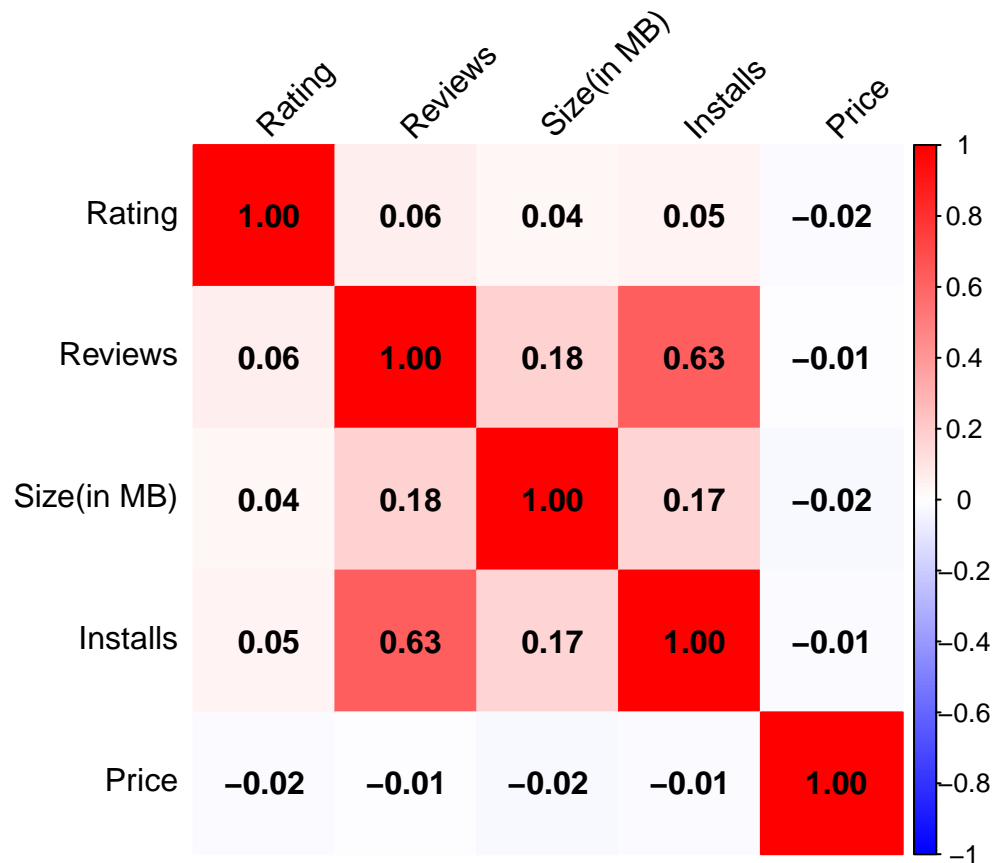
#### Step 4: Distribution of Reviews in Each Category



Similar to installs, I standardize the review variable. As shown in the box plot above, game category still tend to have more reviews than the other category, with the most number of reviews over  $4 * 10^7$ .



## Step 5: Correlation Analysis of Numeric Variables



From the correlation plot, most of the numerical variables do not have strong correlation with each other, except the installs and review. Their correlation is 0.63, which means that there is positive relationship between those two variables. When there is more installs for the app, it tends to have more reviews.

## Step 6: Check Association Between Categorical Variables Using Chi-Square Test

```
##
## Pearson's Chi-squared test
##
## data: Category_vs_Type
## X-squared = 246.46, df = 32, p-value < 2.2e-16

##
## Pearson's Chi-squared test
##
## data: Category_vs_ContentRating
## X-squared = 4377.1, df = 160, p-value < 2.2e-16

##
## Pearson's Chi-squared test
##
## data: Type_vs_ContentRating
## X-squared = 14.532, df = 5, p-value = 0.01256
```

From the chi-square results above, since all the p-value of the tests are less than 0.05, which means that there is association between category and type, category and content rating and type and content rating.

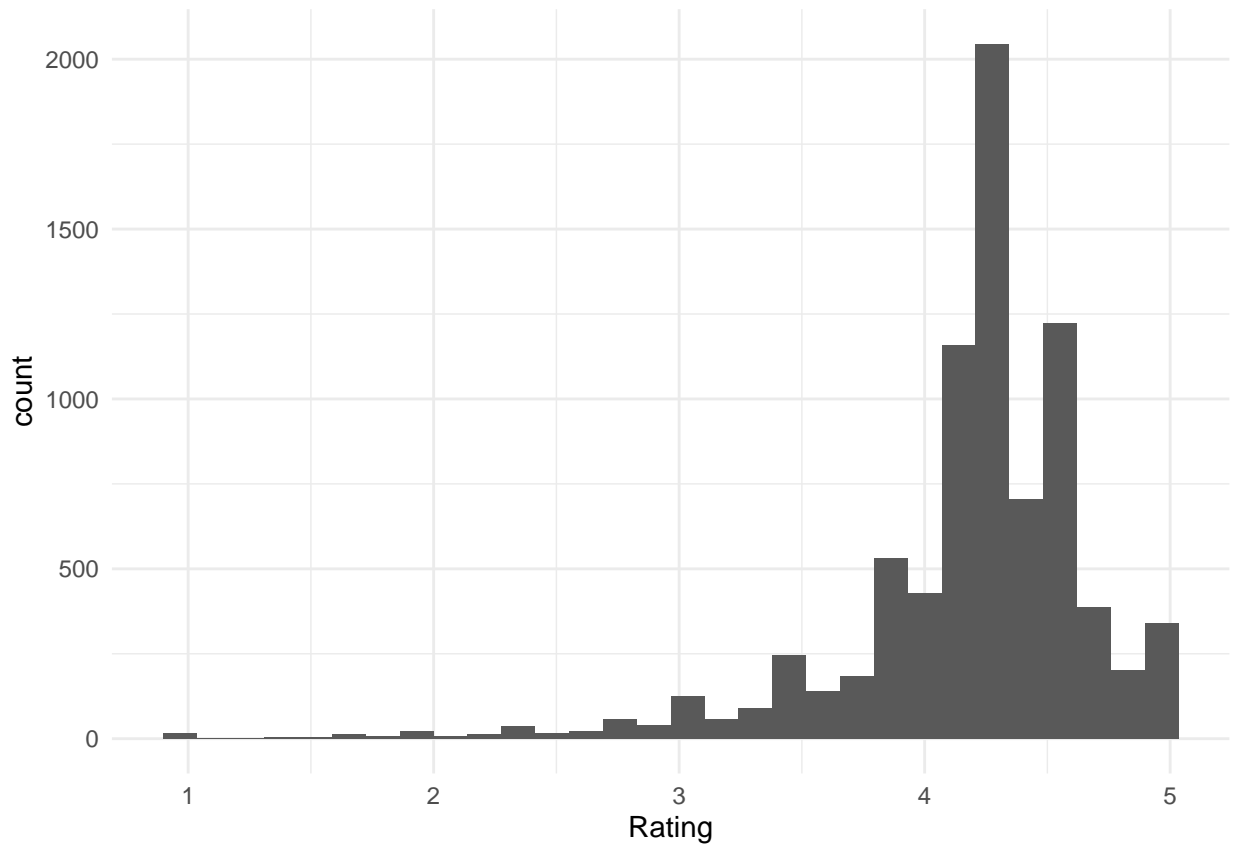
## Method

### Step 1: Selection of Models

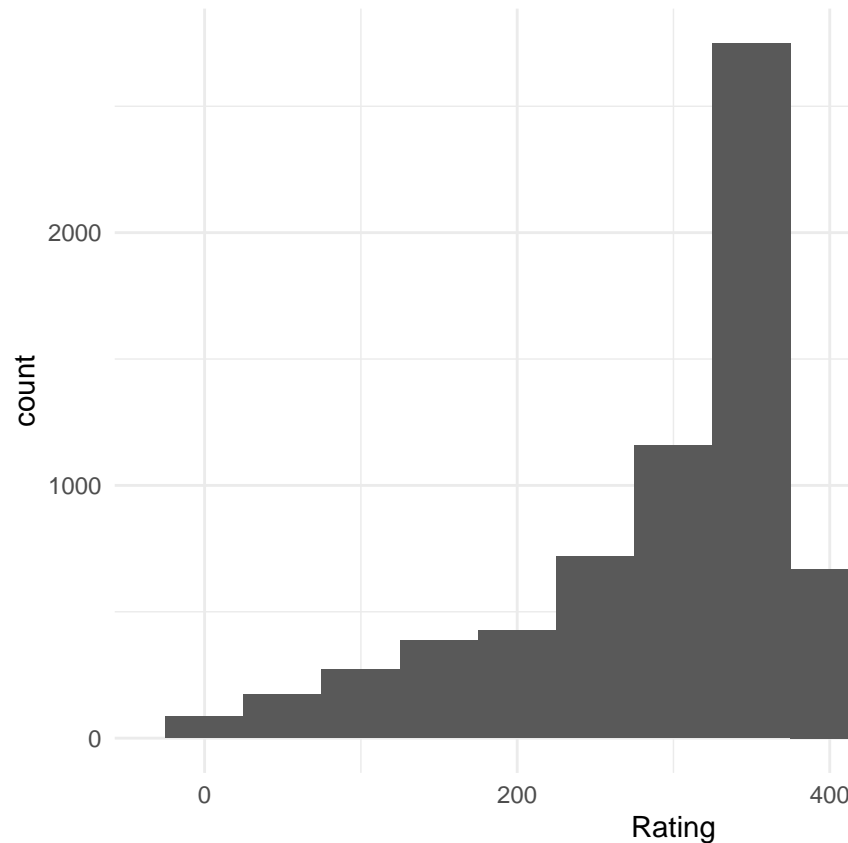
Since the response variable rating is a continuous variable, logistic and multinomial models cannot be used. Also, the response variable is not count data and not discrete, so poisson and negative binomial models cannot be used. Thus, the models that I want to use are linear regression model or linear mixed effect models.

First plot the distribution of rating score:

```
ggplot(data = df, aes(x=Rating))+  
  geom_histogram() +  
  theme_minimal()
```



Since the rating score is left-skewed, I want to do some transformation to make it more symmetric. The



method I would like to use is power transformation.

After taking rating to the power of 4, the distribution of rating score becomes normal. So in the later report, when I refer to the variable rating, it refers to  $rating^4$ .

## Step 2: Start from Null Model

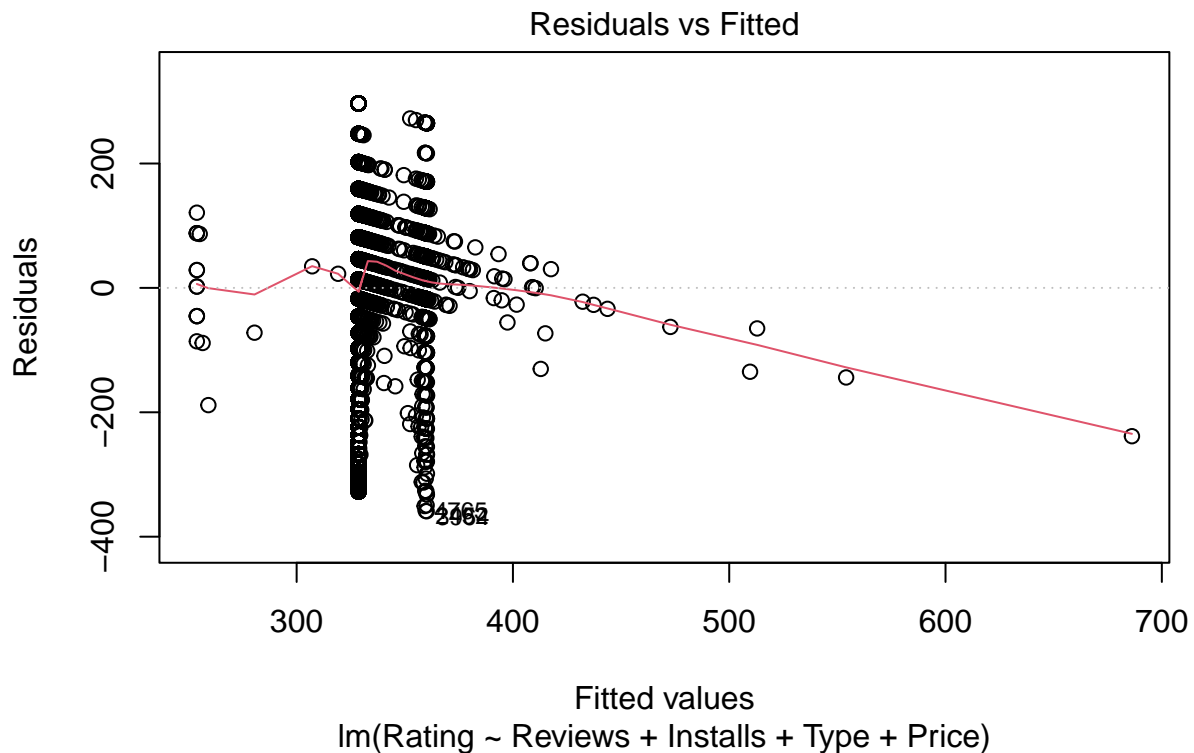
```
##
## Call:
## lm(formula = Rating ~ 1, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -330.77  -75.77   10.11   78.29  293.23
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   331.772      1.358   244.4  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 122.3 on 8117 degrees of freedom
```

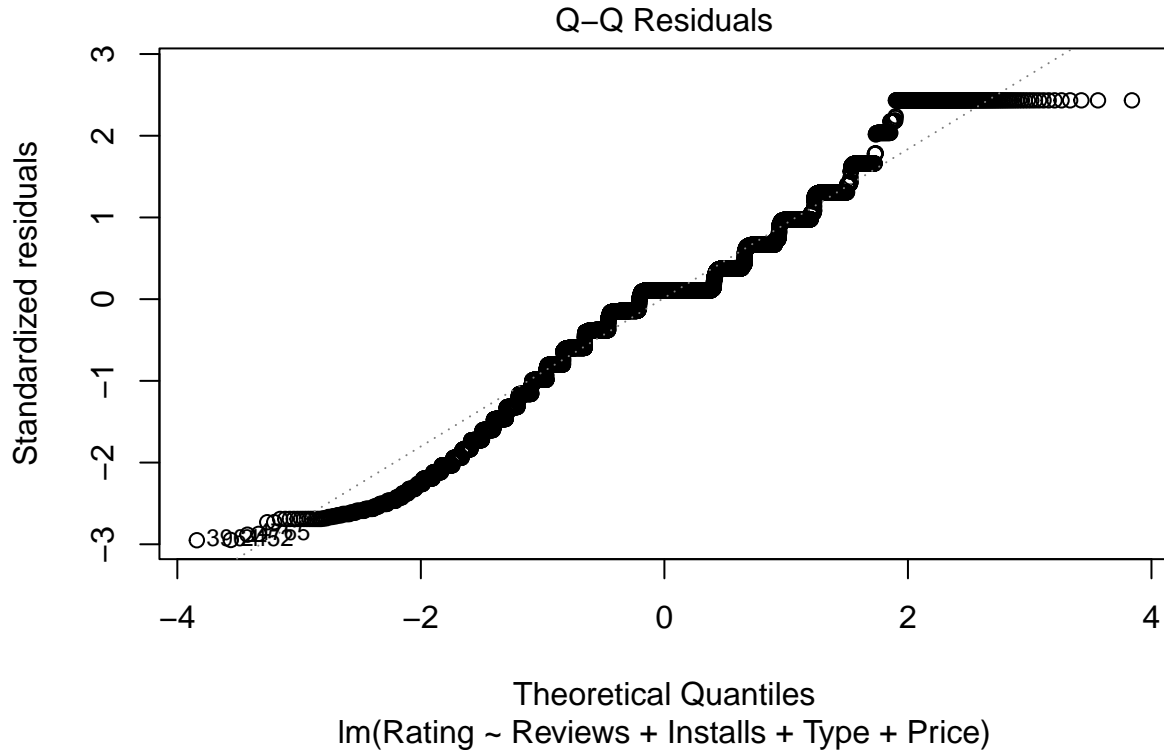
In this null model, the estimate of the intercept is just the overall mean of rating score of all the apps in the whole dataset.

### Step 3: Building the Linear Regression Model and Transformation of Data

Building upon the null model, I will add some variables into the model:

```
##  
## Call:  
## lm(formula = Rating ~ Reviews + Installs + Type + Price, data = df)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -359.10  -72.63   13.25   76.85  296.37   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  3.286e+02  1.434e+00  229.207  < 2e-16 ***  
## Reviews      7.937e-06  1.821e-06   4.358  1.33e-05 ***  
## Installs     5.007e-09  3.677e-08   0.136  0.891707   
## TypePaid     3.187e+01  5.258e+00   6.060  1.42e-09 ***  
## Price       -2.666e-01  7.801e-02  -3.417  0.000635 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 121.8 on 8113 degrees of freedom  
## Multiple R-squared:  0.008725,    Adjusted R-squared:  0.008236   
## F-statistic: 17.85 on 4 and 8113 DF,  p-value: 1.326e-14
```

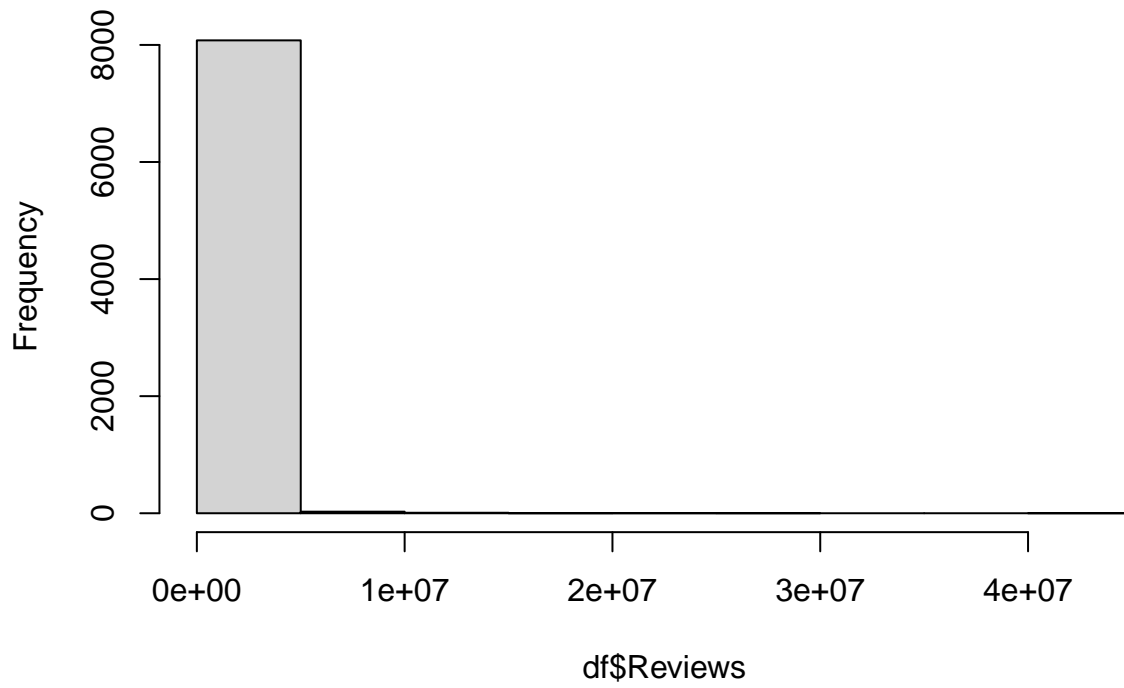


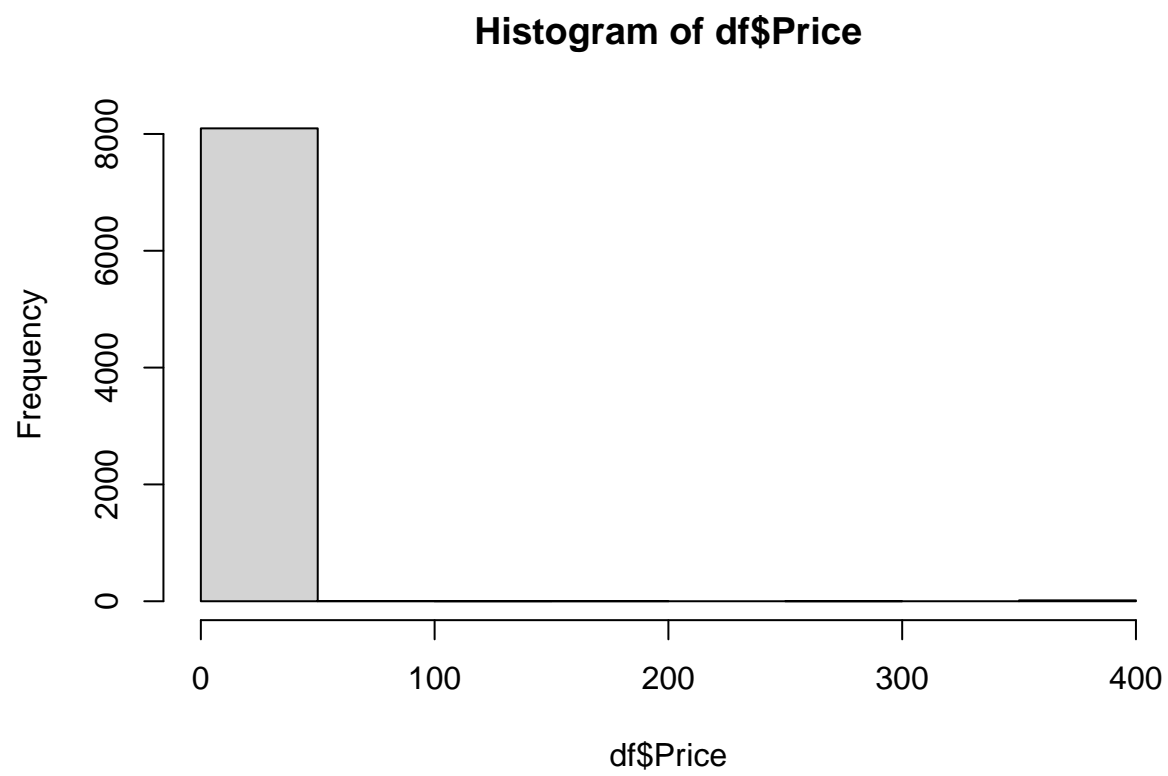


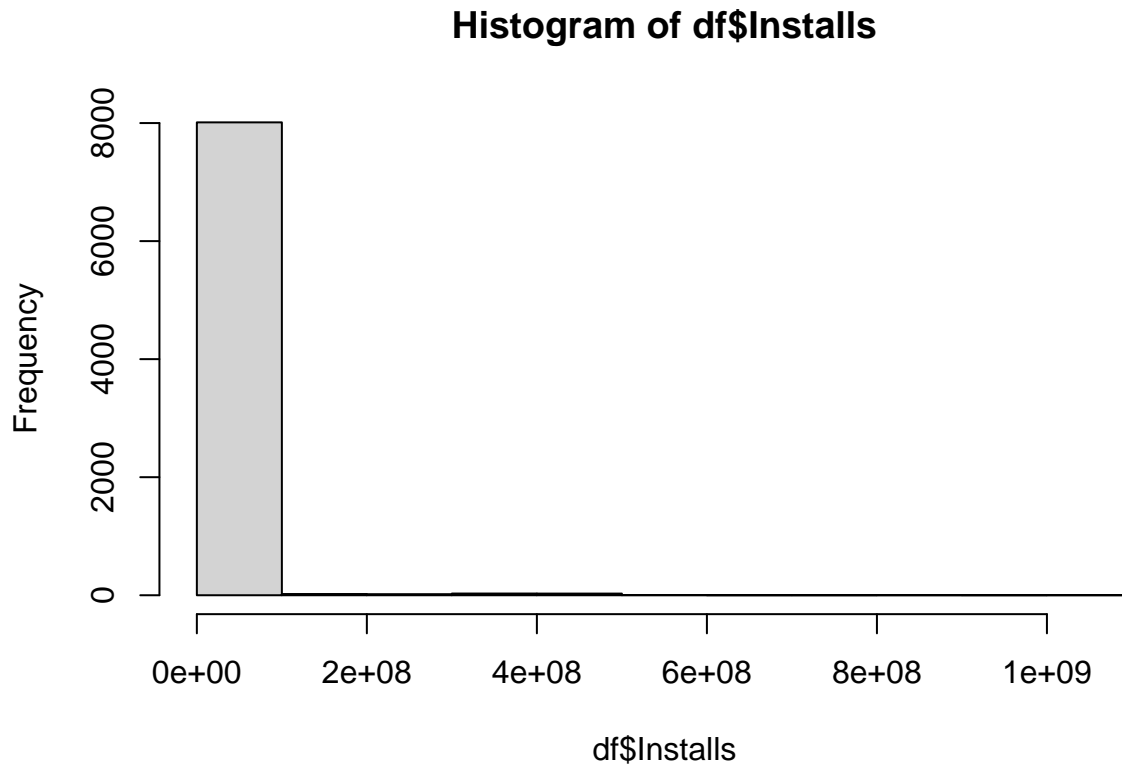
Based on the model result, variable **Reviews**, **Type** and **Price** have p-value less than 0.05, suggesting that they have statistically significant effect on Rating. Additionally, F-test is less than 0.05, suggesting that there is at least one variable that is statistically significant. However, based on the residuals vs fitted value plot, the points are not randomly scattered around 0, indicating that there is heteroscedasticity problem. Based on the QQ plot, there are some points that deviate from the line, but overall, the normality assumption is not violated. Since the homoscedasticity assumptions of the linear regression models is violated, the model cannot be used to fit the dataset, which may explain the low adjusted R-square value.

To tackle this heteroscedasticity problem, one thing I would like to try is doing log transformation on the predictor variable because based on the residual vs fitted plot, the fitted value is concentrated in a narrow range, and this could be related to the high-skewness or low variability in one of the predictor variables:

**Histogram of df\$Reviews**



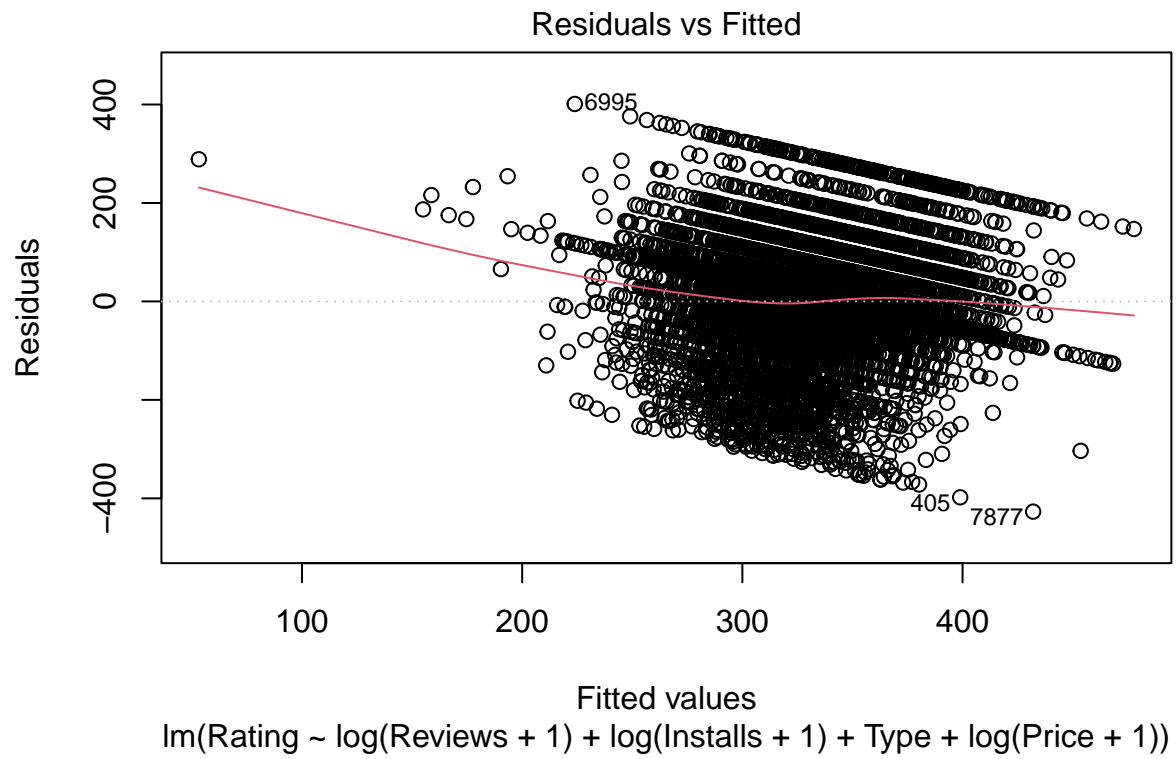


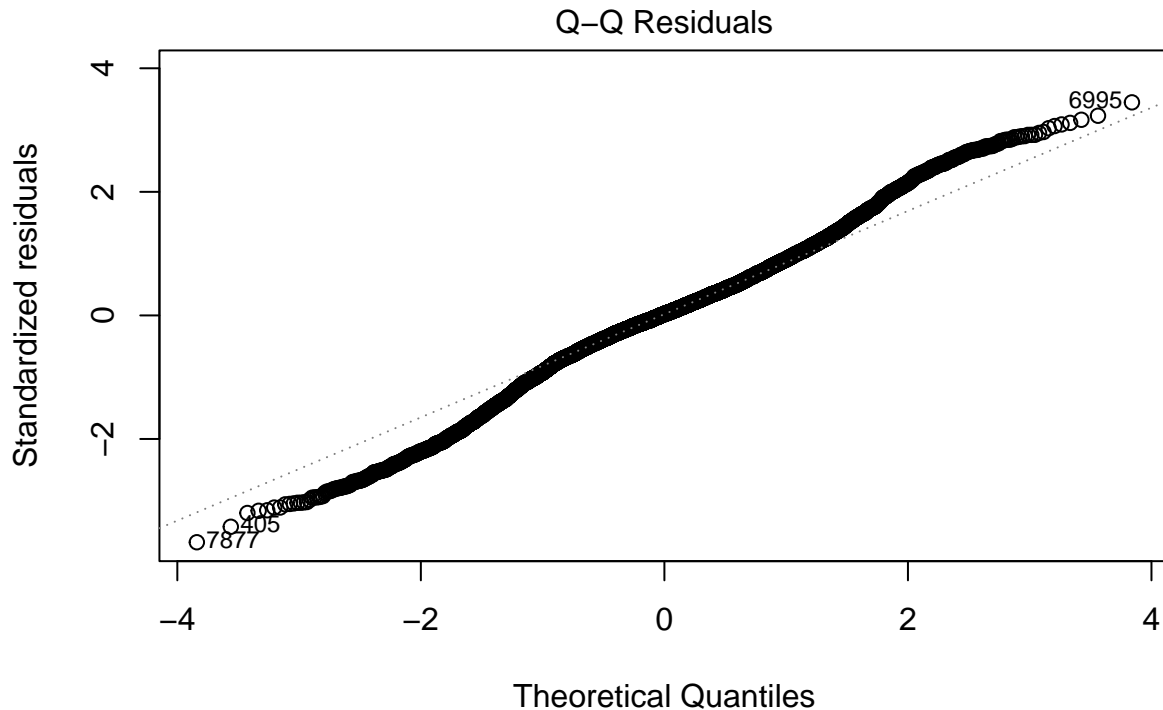


Based on the histograms above, those three numeric variables are all highly skewed, so I could do log transformation on those three variables.

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1), data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -426.98  -63.18    2.68   67.95  401.12
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    467.062     5.637   82.862 < 2e-16 ***
## log(Reviews + 1)  30.666     1.098   27.919 < 2e-16 ***
## log(Installs + 1) -29.840     1.052  -28.352 < 2e-16 ***
## TypePaid         13.204     8.943    1.476 0.139875
## log(Price + 1)   -16.922     4.602   -3.677 0.000237 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 116.4 on 8113 degrees of freedom
## Multiple R-squared:  0.09513,    Adjusted R-squared:  0.09468
## F-statistic: 213.2 on 4 and 8113 DF,  p-value: < 2.2e-16
```







$\text{lm}(\text{Rating} \sim \log(\text{Reviews} + 1) + \log(\text{Installs} + 1) + \text{Type} + \log(\text{Price} + 1))$

After doing the log transformation on predictors, in the residual vs fitted plot, the residuals become more spread around 0, indicating the homoscedasticity assumption is no longer violated, and the QQ plot has less heavy-tail problem, indicating the normality of residual distribution. What's more, the adjust r-squared value is ten times the original model, and the residual standard error also decreased, meaning that model 2 is a better model than model 1. I could add more variables into the model and see if it improves the model, since type is not a statistically significant variable in the model, I will remove it in the later model:

```
model3 <- lm(Rating~log(Reviews+1)+log(Installs+1)+Type+
              log(Price+1)+Category, data = df)
summary(model3)
```

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1) + Category, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -419.18  -63.66    3.69   67.14  404.29
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    532.594    15.871  33.558 < 2e-16 ***
## log(Reviews + 1)    31.905     1.118  28.548 < 2e-16 ***
## log(Installs + 1)  -30.703     1.061 -28.928 < 2e-16 ***
## TypePaid         11.213     8.993   1.247 0.212497
```

```
## log(Price + 1) -17.727 4.605 -3.850 0.000119 ***
## Categoryauto_and_vehicles -61.407 19.873 -3.090 0.002009 **
## [ reached getOption("max.print") -- omitted 31 rows ]
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 115 on 8081 degrees of freedom
## Multiple R-squared: 0.1201, Adjusted R-squared: 0.1162
## F-statistic: 30.64 on 36 and 8081 DF, p-value: < 2.2e-16
```

In model 3, residual standard error decreased and adjusted R-square increased, meaning that adding category variable increase the model performance. Let's try adding two more categorical variables:

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1) + Category + Content.Rating + Last.Updated,
##     data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -375.05  -57.16    0.00   61.54   369.26
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    516.98390   105.36960   4.906 9.49e-07 ***
## log(Reviews + 1)    30.56149    1.21829  25.086 < 2e-16 ***
## log(Installs + 1)  -29.85975    1.13604 -26.284 < 2e-16 ***
## TypePaid          21.14725    10.09366   2.095 0.036199 *
## log(Price + 1)    -17.17115    5.00692  -3.429 0.000608 ***
## Categoryauto_and_vehicles -65.13024   20.26673  -3.214 0.001317 **
## [ reached getOption("max.print") -- omitted 1263 rows ]
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 112.9 on 6849 degrees of freedom
## Multiple R-squared: 0.2806, Adjusted R-squared: 0.1474
## F-statistic: 2.106 on 1268 and 6849 DF, p-value: < 2.2e-16
```

Based on the model4 output result, the residual standard error and adjusted r-squared all get further improvement.

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1) + 'Size(in MB)' + Category + Content.Rating +
##     Last.Updated, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -376.10  -57.01    0.00   61.45   367.58
##
## Coefficients:
```

```
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)      514.64685   105.29008    4.888 1.04e-06 ***
## log(Reviews + 1)    30.98824    1.22375   25.322 < 2e-16 ***
## log(Installs + 1)  -29.99155    1.13582  -26.405 < 2e-16 ***
## TypePaid           22.19363   10.09049    2.199 0.027879 *
## log(Price + 1)     -17.50874    5.00401   -3.499 0.000470 ***
## 'Size(in MB)'      -0.24750    0.07254   -3.412 0.000649 ***
## [ reached getOption("max.print") -- omitted 1264 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 112.9 on 6848 degrees of freedom
## Multiple R-squared:  0.2818, Adjusted R-squared:  0.1487
## F-statistic: 2.117 on 1269 and 6848 DF,  p-value: < 2.2e-16
```

After getting those five models with only fixed effects, I can use MSE to determine which model is the best model and compare it with the null model:

```
## MSE for Null Model: 14959.46

## MSE for Model 1: 14828.94

## MSE for Model 2: 13536.43

## MSE for Model 3: 13162.95

## MSE for Model 4: 10762.43

## MSE for Model 5: 10744.17

## Analysis of Variance Table
##
## Model 1: Rating ~ 1
## Model 2: Rating ~ log(Reviews + 1) + log(Installs + 1) + Type + log(Price +
##           1) + 'Size(in MB)' + Category + Content.Rating + Last.Updated
##   Res.Df      RSS    Df Sum of Sq  Pr(>Chi)
## 1     8117 121440886
## 2     6848  87221189 1269   34219696 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

By comparing their MSE and chi-square test between model 5 and null model, model 5 is the best model, which has the lowest MSE, residual standard error and highest adjusted R-squared.

## Step 4: Fit the Multi-level Model

Based on the model fitting result above, some of the categories of apps are statistically significant, also, in the EDA part, each group's variability in rating score is different. Thus, I could try to treat Category variable as a group variable and fit a LMM on the data.

Let's start with the no pooling model, which is similar to model 5:

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1) + 'Size(in MB)' + Category + Content.Rating +
##     Last.Updated - 1, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -376.10  -57.01    0.00   61.45  367.58
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## log(Reviews + 1)                30.98824     1.22375   25.322 < 2e-16 ***
## log(Installs + 1)             -29.99155     1.13582  -26.405 < 2e-16 ***
## TypeFree                      514.64685    105.29008    4.888 1.04e-06 ***
## TypePaid                      536.84048    105.68525    5.080 3.88e-07 ***
## log(Price + 1)                 -17.50874     5.00401   -3.499 0.000470 ***
## 'Size(in MB)'                  -0.24750     0.07254   -3.412 0.000649 ***
## [ reached getOption("max.print") -- omitted 1264 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 112.9 on 6848 degrees of freedom
## Multiple R-squared:  0.9141, Adjusted R-squared:  0.8981
## F-statistic: 57.36 on 1270 and 6848 DF, p-value: < 2.2e-16
```

In this no pooling model, each category is treated as a separate variable in the model and has its own coefficients. Also, the adjusted R-squared increased a lot compared to model 5, indicating a better fit.

Now let's fit the partial pooling model:

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: Rating ~ log(Reviews + 1) + log(Installs + 1) + Type + log(Price +
##     1) + 'Size(in MB)' + Content.Rating + Last.Updated + (1 | Category)
## Data: df
##
## REML criterion at convergence: 86077.9
##
## Scaled residuals:
##      Min       1Q   Median       3Q      Max
## -3.3263 -0.5045  0.0000  0.5464  3.2204
##
## Random effects:
## Groups Name Variance Std.Dev.
## Category (Intercept) 575.2 23.98
## Residual 12740.9 112.88
## Number of obs: 8118, groups: Category, 33
##
## Fixed effects:
##                                Estimate Std. Error t value
## (Intercept)                  459.33527    104.15285    4.410
## log(Reviews + 1)                30.78869     1.22030   25.230
## log(Installs + 1)             -29.80958     1.13278  -26.315
## TypePaid                      22.67180    10.07540    2.250
```

```
## log(Price + 1)                -17.57672    4.99929   -3.516
## 'Size(in MB)'                 -0.24079    0.07223   -3.334
## Content.RatingEveryone        -38.28120   80.89076   -0.473
## Content.RatingEveryone 10+    -48.25497   81.21448   -0.594
## Content.RatingMature 17+      -59.66625   81.23899   -0.734
## Content.RatingTeen            -51.78637   80.96758   -0.640
## [ reached getOption("max.print") -- omitted 1228 rows ]

##
## Correlation matrix not shown by default, as p = 1238 > 12.
## Use print(x, correlation=TRUE) or
##     vcov(x)           if you need it
```

Based on the output of the partial pooling model, the variance of app rating score of each Category group is very small compared to overall variance of app rating score across all Category groups. Also, compared with no pooling model, the residual standard deviation is similar compared with no pooling model, indicating that the mixed effect may not be necessary.

Another option is to try complete pooling and completely ignore the Category Variable:

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     Type + log(Price + 1) + 'Size(in MB)' + Content.Rating +
##     Last.Updated, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -384.63  -56.48    0.00   61.12  370.96
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    439.16479   105.00472    4.182 2.92e-05 ***
## log(Reviews + 1)    30.43898    1.21390   25.075 < 2e-16 ***
## log(Installs + 1)  -29.53540    1.13105  -26.113 < 2e-16 ***
## TypePaid         27.49832    10.03912    2.739 0.006176 **
## log(Price + 1)    -18.58360    5.00152   -3.716 0.000204 ***
## 'Size(in MB)'     -0.13095    0.06775   -1.933 0.053307 .
## [ reached getOption("max.print") -- omitted 1232 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 114.3 on 6880 degrees of freedom
## Multiple R-squared:  0.2604, Adjusted R-squared:  0.1275
## F-statistic: 1.958 on 1237 and 6880 DF, p-value: < 2.2e-16
```

In this model, compared with previous two models, the residual standard error increased, and its adjusted R-squared is lower than the no pooling model, indicating that the model may be worse than the previous two.

To further compare those three model, MSE could be calculated:

```
## MSE (No Pooling): 10744.17
```

```
## MSE (Partial Pooling): 10756.35
```

```
## MSE (Complete Pooling): 11063.66
```

Based on the MSE calculated, the model with no pooling gets the lowest MSE, which means that it is the best model among the three models.

## Step 5: Compare the No Pooling Model and Model 5

To compare the performance of those two models, we can first

```
## AIC (model5): 100932.1
```

```
## MSE (model5): 10744.17
```

```
## AIC (No Pooling): 100932.1
```

```
## MSE (No Pooling): 10744.17
```

Based on the results above, no pooling model gets very similar result with model 5 for both AIC and MSE because no pooling model is just model 5 removing the intercept term. For easier model interpretation, I would choose model 5 for further analysis, but with high AIC value and lower MSE value, this could indicate the overfitting problem with model 5, but with the knowledge in MA678, I don't have a way to verify this.

## Step 6: Add Interaction Term into No Pooling Model

Since based on the model 5 output, type and size are statistically significant, I could try adding the interaction term between these two:

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##     log(Price + 1) + 'Size(in MB)' * Type + Category + Content.Rating +
##     Last.Updated, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -376.29  -57.05    0.00   61.14  366.81
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    515.16570   105.21837    4.896 1.00e-06 ***
## log(Reviews + 1)    31.27185    1.22609   25.505 < 2e-16 ***
## log(Installs + 1)  -30.21603    1.13719  -26.571 < 2e-16 ***
## log(Price + 1)    -16.36018    5.01332   -3.263 0.001106 **
## 'Size(in MB)'    -0.31060    0.07510   -4.136 3.58e-05 ***
## TypePaid          3.09937    11.70020    0.265 0.791096
## [ reached getOption("max.print") -- omitted 1265 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 112.8 on 6847 degrees of freedom
## Multiple R-squared:  0.2829, Adjusted R-squared:  0.1498
## F-statistic: 2.127 on 1270 and 6847 DF,  p-value: < 2.2e-16

##      Estimate Std. Error    t value    Pr(>|t|)
## 0.761926630 0.236797722 3.217626521 0.001298578

## Analysis of Variance Table
##
## Model 1: Rating ~ log(Reviews + 1) + log(Installs + 1) + log(Price + 1) +
##      'Size(in MB)' * Type + Category + Content.Rating + Last.Updated
## Model 2: Rating ~ log(Reviews + 1) + log(Installs + 1) + Type + log(Price +
##      1) + 'Size(in MB)' + Category + Content.Rating + Last.Updated
##   Res.Df      RSS Df Sum of Sq Pr(>Chi)
## 1    6847 87089504
## 2    6848 87221189 -1    -131685 0.001293 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Based on the model output result and anova table result above, although the interaction term is statistically significant, and adding the interaction term does improve the model's adjusted R-square a little bit, the chi-square test still show that model 2, which is the model without interaction term, is better. Thus, the additional interaction term between Size and Type is not needed.

## Step 7: Interpretation of Model 5 Selected

Now, after choosing the best model, I will try to interpret this model:

```
##
## Call:
## lm(formula = Rating ~ log(Reviews + 1) + log(Installs + 1) +
##      Type + log(Price + 1) + 'Size(in MB)' + Category + Content.Rating +
##      Last.Updated, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -376.10  -57.01    0.00   61.45   367.58
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    514.64685   105.29008   4.888 1.04e-06 ***
## log(Reviews + 1)    30.98824    1.22375  25.322 < 2e-16 ***
## log(Installs + 1)  -29.99155    1.13582 -26.405 < 2e-16 ***
## TypePaid          22.19363   10.09049   2.199 0.027879 *
## log(Price + 1)    -17.50874    5.00401  -3.499 0.000470 ***
## 'Size(in MB)'     -0.24750    0.07254  -3.412 0.000649 ***
## [ reached getOption("max.print") -- omitted 1264 rows ]
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 112.9 on 6848 degrees of freedom
```



```
## Multiple R-squared:  0.2818, Adjusted R-squared:  0.1487
## F-statistic: 2.117 on 1269 and 6848 DF,  p-value: < 2.2e-16
```

Based on the model output, the residual's median is 0, and min, max, first quantile and third quantile are all close to each other, meaning that the residual's distribution is symmetric.

To interpret the coefficient, since there are many variables, especially in category, content rating and last updated, I will only select one of the level from each variable to interpret.

For the coefficient of reviews, one unit increase in  $\log(\text{Reviews} + 1)$ , the ratings(which is the 4th power of original rating score) is expected to increase by 31.184, keeping all other variables constant. This shows the positive correlation between reviews and rating score.

For the coefficient of installs, one unit increase in  $\log(\text{Installs} + 1)$ , the ratings is expected to decrease by 30.14, keeping all other variables constant. This shows the negative correlation between installs and rating score.

The interpretation of coefficient of price and size are similar to previous ones, and they both have negative correlation with the rating scores.

For coefficient of Type, the average rating score for paid app is 21.32 higher compared to free app, keeping all other variables constant.

For category level auto and vehicles, the coefficient is -58.2, which shows that the average rating score for apps belonging to auto and vehicles group is 58.2 lower compared to the apps in arts and design group(reference group), keeping all other variables constant.

For content rating level everyone, the coefficient is -42.12, which shows that the average rating score for apps with content rating as everyone is 42.12 lower compared to the apps with content rating as Adults(reference group), keeping all other variables constant.

For the last updated date April 1st, 2017, the coefficient is 54.42, which shows that the average rating score for apps that last updated at April 1st, 2017 is 54.42 higher compared to the apps last updated at April 1st 2016(reference group), keeping all other variables constant.

The residual standard error is 112.8, which is the lowest compared to other models that I fit. The adjust R-square is 0.1492, which indicated that the model can only explain 14.92% variability in the response variable, and although this is the best model I get, the model still does not fit the data well. The F-statistic has p-value less than 0.05, showing that there is at least one variable that is statistically significant, and the model is better than the null model.

The important predictors of rating score according to my model are reviews, installs, type, price, size, category and last updated date.

## Discussion

To validate the result of my model, I try to find three research paper that related to rating score prediction in Google Play Store.

Based on the paper by S Shashank et.al[2], they tried to predict the rating score using machine learning algorithm. The paper used the same dataset as I do in this report. Compared to my approach, the author do a more detailed EDA. Instead of focusing on the category of apps, they pay more attention on whether the app is free or paid and find out the difference in rating score between free app group and paid app group. In the method part, the authors applied five techniques trying to find the important variables related to rating score: random forest, support vector regression, linear regression, k-nearest neighbors and k-means clustering. As a results, the k-means neighbors achieves the best result, which 92% prediction accuracy, and the author concludes that **size, type, price, content rating and genre** are variables strongly correlated with the rating score. The author's result is consistent with my result, and those are also statistically

significant variables in my linear regression model, except I used category instead of genre. However, after reading the paper, I find out there could be bias related to the rating score because higher ratings given by users potentially attract several new users disproportionately, and people tend to only use apps with high rating score, leading to more reviews of the app. Additionally, many people don't like writing reviews for the apps no matter they like the app or not, there are also people writing negative reviews but give very positive rating scores, so for some apps, the rating score may not reflect their true quality.

Another research paper I found was written by Min-Kyo Seo et.al[3]. The main purpose of this paper is to investigate the predictors and main determinants of consumers' ratings of mobile applications in the Google Play Store. The author also tried to extend their model into a sentimental analysis and aim to review polarity and subjectivity on the application rating. In the data preprocessing part, they used sentiment analysis based on the users reviews from Google Play Store and created new variables polarity and subjectivity and merge them into the original Google Play dataset. In the method part, there were four models they used: multiple linear regression, regression tree, random forest tree and neural network. Based on the model result, neural network model gives the lowest RMSE result, and the important variables are **price, installs and reviews**, and polarity and subjectivity of reviews are less critical. Those variables are also included in my model, one difference is that install is positively correlated with ratings based on the author's results, but install is negative correlated with rating score based on my model's results.

The last research paper I found is by Jayanth. P et.al[4], and the main purpose of research paper is to predict the rating score using the comprehensive Google Play Store dataset similar to the one that I used. The method that the author used are: lasso regression model, ridge regression model, gradient boosting, XGBoost and CATBoost. Based on the result, CATBoost method provides the lowest MAE, MSE, RMSE and highest  $R^2$  value, but the top features that the author find that can influence the rating score is different from my results, the author find that **reviews, last updated and android version** are the most critical variables. In my models, android version is not included, and most of the levels of last updated are not statistically significant.

Based on the literature review above, the next step for my analysis could be try to implement tree model or deep learning model to further improve my current linear regression model.

## Appendix

### Reference

- [1] Statista (2024) Google Play Store - Statistics & Facts. Available at: <https://www.statista.com/topics/9929/google-play-store/#topicOverview> (Accessed: 27 November 2024).
- [2] S Shashank and Brahma Naidu, "Google play store apps-data analysis and ratings prediction", International Research Journal of Engineering and Technology, vol. 7.12, pp. 265-274, 2020.
- [3] Seo, M.K., Yang, O.S. and Yang, Y.H., 2020. Global Big Data Analysis Exploring the Determinants of Application Ratings: Evidence from the Google Play Store. Journal of Korea Trade, 24(7), pp.1-28.
- [4] J. P. A. Nagam, P. Undavalli, P. P. V. P. K. S and V. K. K. K, "Leveraging CAT Boost for Enhanced Prediction of App Ratings in the Google Play Store," 2024 Second International Conference on Advances in Information Technology (ICAIT), Chikkamagaluru, Karnataka, India, 2024, pp. 1-6, doi: 10.1109/ICAIT61638.2024.10690600.