

Statistical report

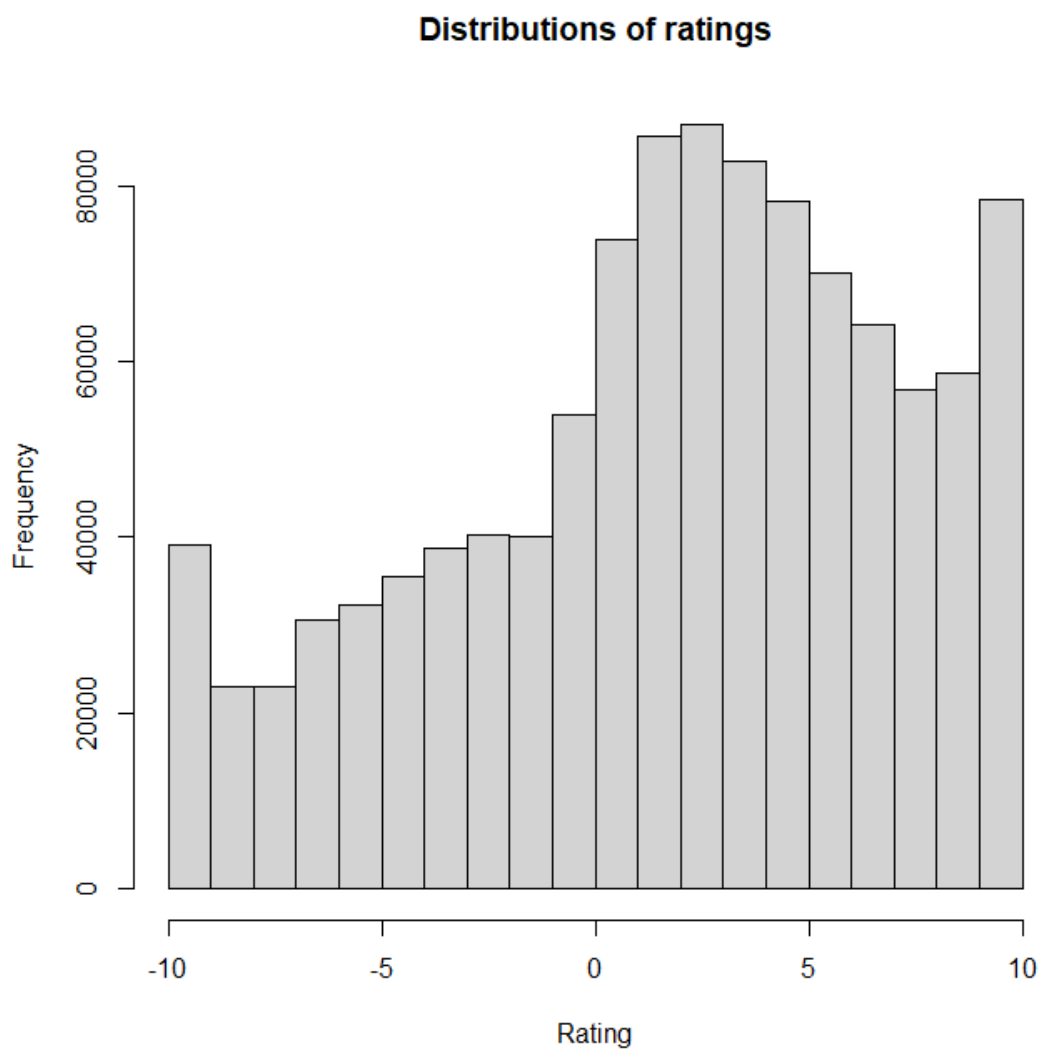
This document is a statistical report on the dataset used in our project. The dataset contains anonymous ratings(-10 to 10) provided by a total of 41,000 users. Train file contains 1.1 million ratings for 139 jokes.

Here we load the data and create the dataframe.

```
library("ggplot2")  
trainSet <- read.csv("train.csv")  
df <- data.frame(trainSet)
```

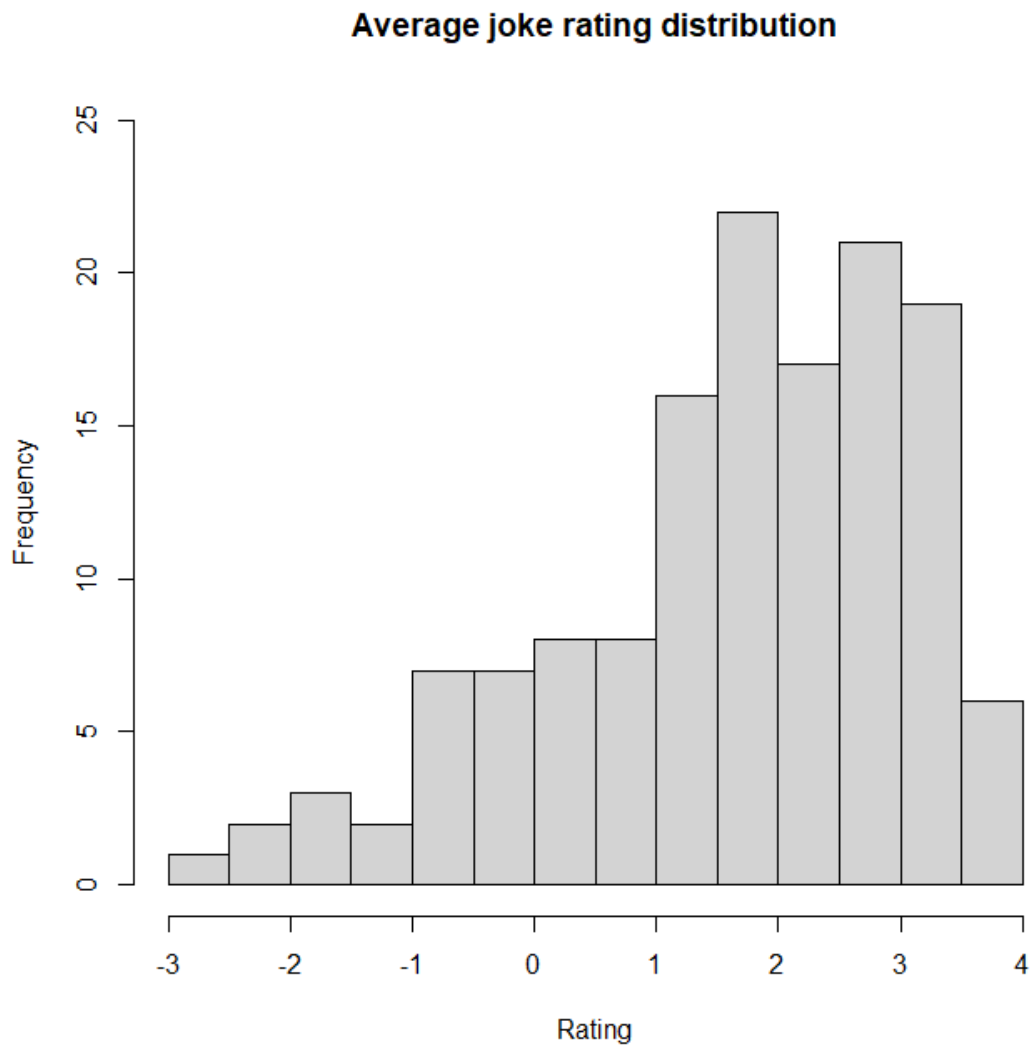
Histogram below represents distribution of the grade. The distribution is little skewed to the right.

```
hist(df$Rating,main = "Distributions of ratings", xlab = "Rating")
```



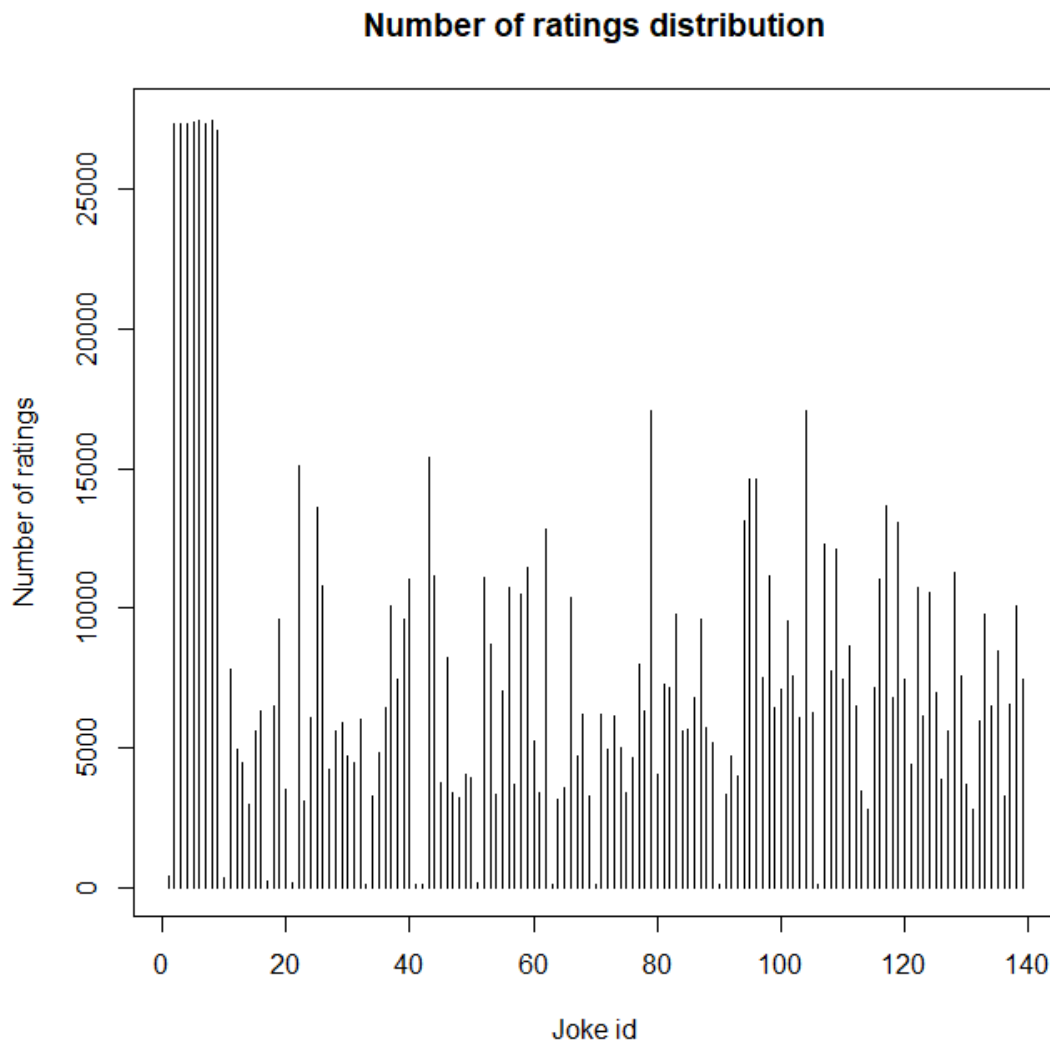
Histogram below represents average joke rating distribution. Average joke rating ranges from -3 to 4, the most popular are within [1.5, 2]. The distribution of average ratings is a little skewed to the right.

```
average_ratings <- aggregate(df$Rating ~ df$joke_id, df, mean)
hist(average_ratings$`df$Rating`, ylim=c(0,25), breaks = 20, main = "Average joke
rating distribution", xlab = "Rating")
```



Histogram below represents number of ratings distribution. As we can see, the popularity of jokes is different. While some jokes have more than 25,000 ratings, others do not even reach 1000 ratings.

```
ratings_numbers <- aggregate(df$Rating ~ df$joke_id, df, length)
plot(ratings_numbers, type = "h", breaks = 139, main = "Number of ratings
distribution", xlab = "Joke id", ylab = "Number of ratings")
```



Histogram below represents standart deviations of jokes ratings distribution. As we can see, the standart deviations form a normal distribution. From this we can conclude that the ratings of jokes are compatible.

```
standart_deviations <- aggregate(df$Rating ~ df$joke_id, df, sd)
hist(standart_deviations$`df$Rating`, xlim = c(4, 6.5), ylim=c(0,20), breaks =
20, main="Standart deviations distribution", xlab = "Standart deviation",
ylab = "Jokes number")
```

