

Datascience Fundamentals with R - Reddit

Group 7

December 5, 2017

Goal

Reddit is a famous American social news aggregation and discussion website with a community of 1.6 billion users. The members post content on the website and it is then voted up or down by other members of the community. This concept is very interesting since it gives a lot of information on what people are interested in and what matters for them at a particular point in time. Thus, it is the reason why we chose to work with Reddit data.

For the project, our group is particularly interested in the reach of a comment defined by its number of upvotes and comments. A second pillar of the project is based on Natural Language Processing based on the post content in order to extract meaningful information from the users' behavior and from the posts content.

In this paper, we first start to analyze the data set by showing the distribution of important post features, we then look at how the number of comments influences the upvotes and how the time passed since posting time impacts the change in number of upvotes.

Then, we put emphasis on Natural Language Processing in order to analyze the users' behavior and the posts. We first describe the sentiment of the data and show how the authors tend to repost news posts on Reddit. We then estimate the proportion of each sentiment in reddit titles and how the posts' sentiments are related to each other.

In a third part, we apply the perceptron model to classify successful posts based on posts' features. Finally, The output of the paper gives a great insight in the drivers of the reach of a post and how people post on Reddit.

Data

We decided to use one month of Reddit News data. We picked December 2016 for this as it was just after the US elections and that enough time has passed for us to (possibly) notice the effects of news that came out at that time.

Reddit news data usually comes in the form of a title and a link, the title is written by the user and links to the news website's page. It sometimes has thumbnails.

Reddit post data is available on google's bigquery database, it is saved in `news_2016_12.csv`. Following are the most relevant columns for our analysis:

- **time_created** (UTC timestamp) - when was the post created
- **author** - username of user that posted
- **domain** - which domain did the news come from?
- **url** - Specific URL of the news post
- **score** : upvotes - downvotes (renamed as **like_score**)
- **upvotes** : how many "likes" the post received
- **downvotes** : howmany "dislikes" the post received
- **title** : user-created title of the post

Due to an issue with the API, no downvote data is given to us. Thus the **score** is equal to the **upvotes** of a post. We will be using the **score** field moving forward.

Furthermore, we thought it would be interesting to compare the difference between the user-generated title and the Actual title posted by the news agency. For this, a (scrapy)[<https://scrapy.org/>] spider was created to crawl all the (cleaned) URLs and retrieve the title. We think this was rather successful since it retrieved 24,045 titles from about 31,713 cleaned posts. This is saved as `titles.csv`.

Task

Structure of the project

1. Cleaning (in common with the other group)
2. Descriptive analysis of the data set: Exploring the data
3. Introduction of sentiment analysis
4. Prediction of reach with Perceptron
5. Convolutional network to predict subreddit classification

Due to the size and variety of elements in our dataset, we will apply a variety of methods to extract information, and may use external data to improve our analysis. However, this also implies that the performance measures of our algorithms will vary with the specific relationship under examination.

1. Cleaning [1_clean.r]

Everything related to cleaning the original dataset is defined in `clean.r`. Essentially, the script:

- saves `utc_created` as a POSIXct variable
- removes special characters from titles
- removes non-english titles, which halves our original dataset.

Getting relevant data - [05_Classification_Data.R]:

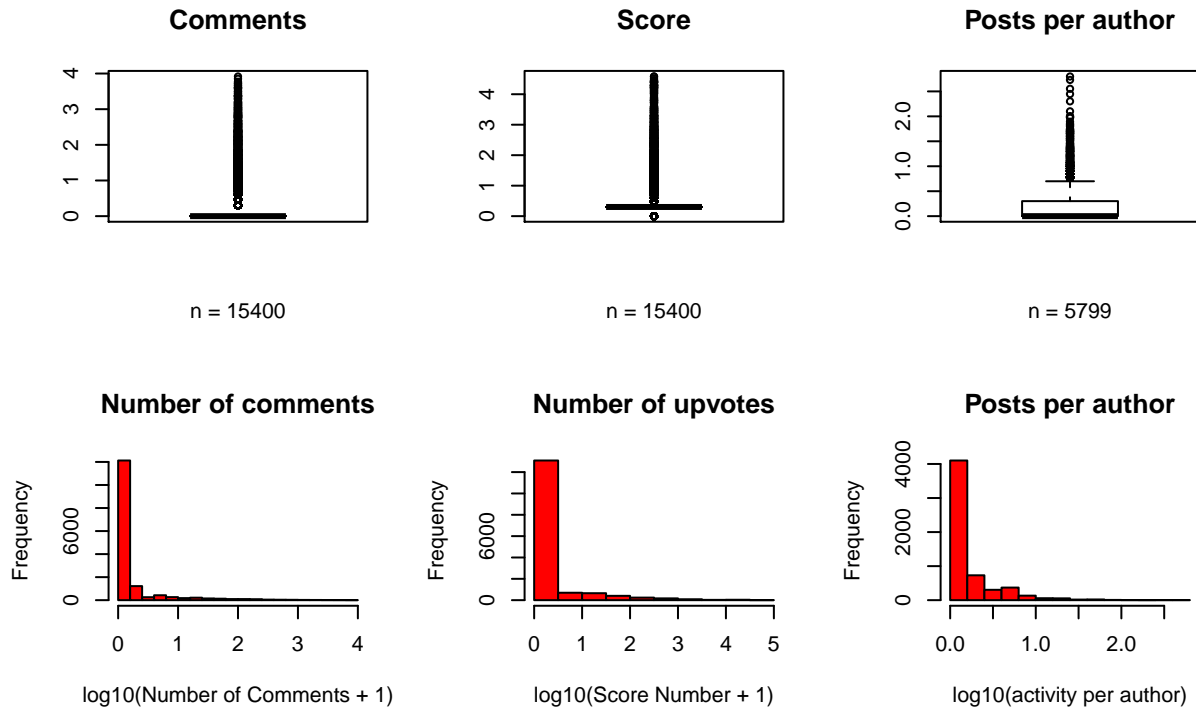
We remove all authors with the value `[deleted]`, because this is value given to users who have deleted their accounts, which mean that we become unable to different specific users.

After these two operations, our dataset is reduced to 3946 rows compared to the original 31,713, and this is further reduced to 394 upon generating the `author_score`

2. Descriptive analysis of the data set: Exploring the data [2_analysis.r]

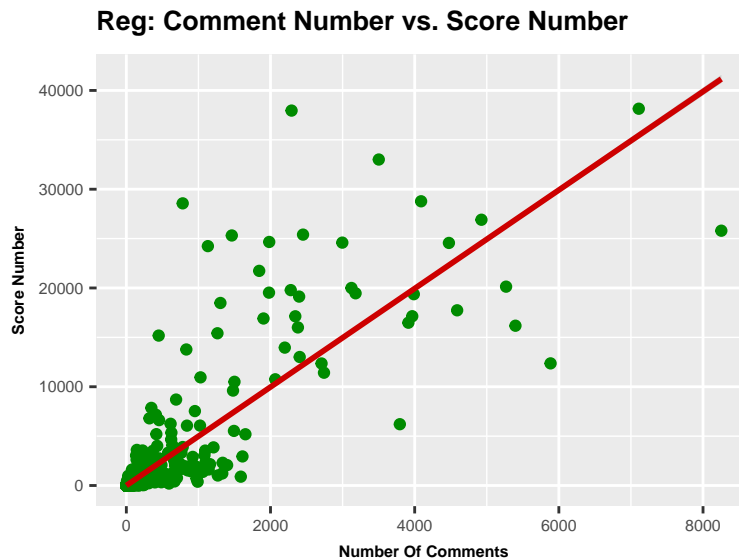
2.1 Analysis of comments, upvotes and posts per author:

Since the data is massively skewed, we set the proportion as \log_{10} .



Interpretation: the plots show a massive skew even though the data's proportion is log10. All plotted variables are right side skewed which shows that the sample consists many low discrete values and marginal number of outliers. This is very important since it has to be taken into consideration in order to adapt the data to get significant results later in this project depending on the models.

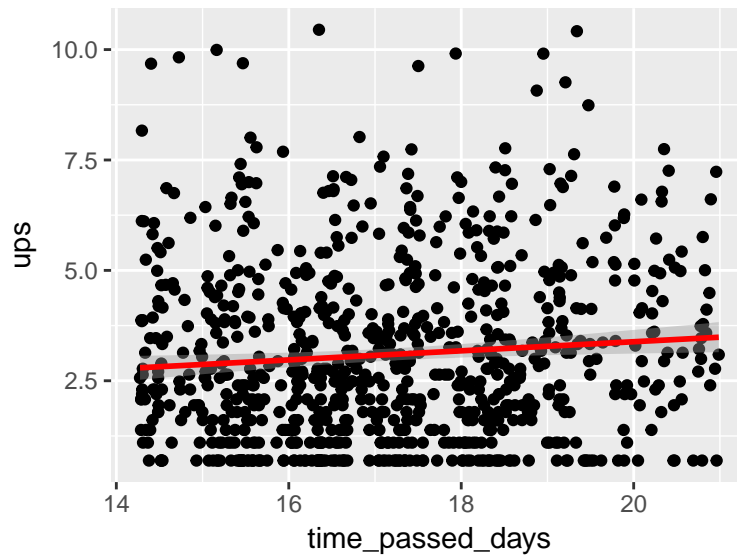
2.2 Regression: Analyze how activity, here defined as number of comments, influences score.



Interpretation: Here, we can see that an additional comment yields on average a 4.9 increase in the score number. The result is statistically significant given the low p-value of $< 2.2e-16$.

2.3 Regression: Analyze how the time passed influences the number of upvotes.

Explanations: here, we log the number of upvotes because the distance between the different number of upvotes is too high and we don't get any significant result. We also select only the posts which have more than one upvotes and focus on posts which are 7 days old because the upvotes usually happens in the next days following the posting date.



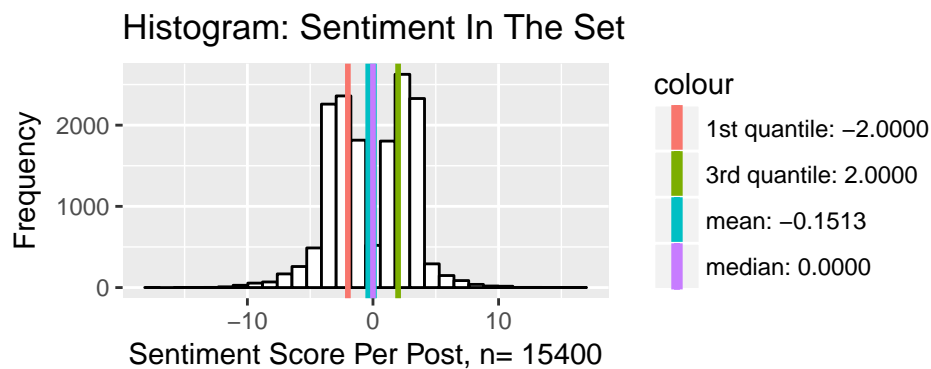
Interpretation: For every additional day, a post gets an additional 0.09013 upvote. The result is statistically significant given the p-value < 0.02902 .

3. Sentiment analysis

Natural Language Processing with tidytext:

In this part, we will use 2 different libraries: NRC and Afinn in order to compute the sentiment of each word of the data set and ultimately be able to compute the sentiment of the different titles. These 2 libraries will be used for two very different purposes, which will be illustrated in the following examples.

3.1 Histogram of sentiment analysis of the data set using the Afinn dictionary



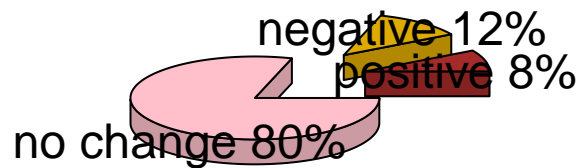
Interpretation:

3.2 Analyze posting behavior with sentiment analysis using AFINN:

In this part, the goal is to analyze whether people interpret news in a positive or negative way when reposting a news. This is done by analyzing the difference of sentiment level between the source title (title displayed on a news website such as on times.com) and the actual title of the reddit user's post.

Approach: to do this, source titles have been retrieved with the help of their URL, then the sentiment score on both title lists has been calculated with the help of the AFINN lexicon assigning a weight between $-5 < \text{weight} < 5$ for every word. We chose AFINN because it ranks the sentiment on a scale, making it easier to compare. Finally, the difference has been computed which enables to get some insights as seen in the following graphics: the histogram shows the proportions of posts unchanged, positive or negative while the 3D chart illustrates the percentage of these classes.

Pie Chart 3D of interpretation levels

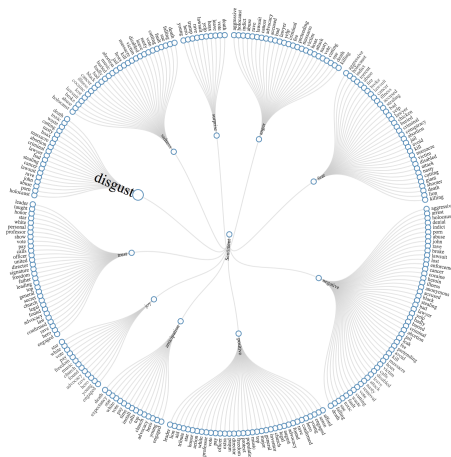


Interpretation: As observed on the pie chart, most Reddit users (4 out of 5 users) tend to repost exactly the same title that comes from the source title, while 12% of the users in this data set repost it in a negative way (when different in sentiment is negative) and the rest (9%) repost it in a positive way.

2.3 Estimation of proportion of each sentiment in reddit titles: how to illustrate sentiment through : Radial dendrogram

There are many ways to represent and analyse the data sets, we decided to explore the beauty of our Reddit dataset in several ways. At first, we used a 'visNetwork' library to create a node-based radial dendrogram. It requires creation of the hierarchy, so we modified our initial sentiment dataset with additional column containing the following hierarchy: Sentiment | Sentiment Detail | Word. The following diagram shows the sentiment allocation for the words used in top 50 posts.

It's interactive. Click [here](#) to try it out online.



2.4. How the sentiments are related : network diagraph

```
# Tatiana's dendrogram
```

4. Perceptron classification:

4.1. Machine learning with Perceptron: prediction of a successful post based on 3 features

Feature 1: sentiment of the post (afinn dictionary) Feature 2: length of the post (number of characters)
Feature 3: Posting time (within a day)

In this part, the three features illustrated above will be used to predict whether a post is successful or not. First, we define a successful post as a post getting at least 2 likes.

This splits the data in X% successful posts and X% not successful posts. Then, we prepare the data by assigning every successful post the category 1 and every non successful post the category -1. After this, we keep only the posting time in hours as well as the sentiment of each title. The Training data is then split in 70% training and 30% testing.

Let's have a look at the training

```
#show the learning effect
```

Performance of the algorithm

```
##summary of results for non classified
```

Discussion on the implementation