

Homework 6: Prediction of rating from Yelp review text

Grover, Karan & Arora, Pragya & Ghai, Piyush

{grover.120, arora.170, ghai.8}@osu.edu

November 29, 2016

1 PROBLEM STATEMENT

1.1 ABOUT YELP

Yelp [1] a famous website as well as a mobile app which publishes crowd sourced reviews about food joints and businesses. It also has a division which handles online reservations for restaurants. **Yelp Dataset Challenge**[2] is a publicly open contest sponsored by Yelp, in which the participants are challenged to use Yelp's data in an innovative way.

1.2 OUR MISSION

The Yelp dataset downloaded from Yelp dataset challenge website is huge and consists of over 2.7M reviews by roughly 687k users for over 86k businesses [2]. For this project, we chose to predict a review's rating based on the review text. The rating will be done on a scale of 1-5, where 1 stands for awful and 5 stands for excellent. We built multiple models and accessed which models would fit our use case the best. This is explained in more depth in the later sections of this report.

1.3 ABOUT THE DATASET

The Yelp Dataset consists of several files in JSON format of the data. The main files as per our use-case were the *yelp_academic_dataset_business.json* & *yelp_academic_dataset_review.json*. The data representation in the for *academic_dataset_business* is as follows :

```
1 {
2     "type": "business",
3     "business_id": (encrypted business id),
4     "name": (business name),
5     "neighborhoods": [(hood names)],
6     "full_address": (localized address),
7     "city": (city),
8     "state": (state),
9     "latitude": latitude,
10    "longitude": longitude,
11    "stars": (star rating, rounded to half-stars),
12    "review_count": review count,
13    "categories": [(localized category names)]
14    "open": True / False (corresponds to closed, not business hours),
15    "hours": {
16        (day_of_week): {
17            "open": (HH:MM),
18            "close": (HH:MM)
19        },
20        ...
21    },
22    "attributes": {
23        (attribute_name): (attribute_value),
24    },
25 }
```

The data representation in the for *academic_dataset_review* is as follows :

```
1 {
```

```
2     "type": "review",
3     "business_id": (encrypted business id),
4     "user_id": (encrypted user id),
5     "stars": (star rating, rounded to half-stars),
6     "text": (review text),
7     "date": (date, formatted like "2012-03-14"),
8     "votes": {(vote type): (count)},
9 }
```

The given datasets were first converted and exported into csv formats using a simple python script. The python script is a part of the **PreProcessing1.py** python file.

REFERENCES

- [1] ggplot Library is used in this assignment to plot most of the graphs in this assgnment <http://ggplot2.org/>.
- [2] ggplot Library is used in this assignment to plot most of the graphs in this assgnment <http://ggplot2.org/>.