

Contents

| | | |
|----------|--|-----------|
| 1 | Introduction | 3 |
| 1.1 | Problem description | 3 |
| 1.2 | Datasets | 3 |
| 1.3 | Tools used | 3 |
| 2 | Preprocessing | 5 |
| 2.1 | Data cleaning | 5 |
| 2.2 | Fake Review Detection | 5 |
| 2.3 | Historical features | 7 |
| 2.4 | User-based collaborative approach | 7 |
| 2.5 | Dimensionality reduction and further preprocessing | 8 |
| 3 | Models | 11 |
| 3.1 | Linear Support Vector Machine | 11 |
| 3.2 | Random Fores | 11 |
| 3.3 | CNN??? | 11 |
| 4 | Experiments and evaluation | 12 |
| 4.1 | Linear Support Vector Machine | 12 |
| 4.2 | Random Forests | 12 |
| 4.3 | CNN??? | 12 |
| 5 | Conclusions | 13 |
| 6 | References | 14 |

List of Tables

| | | |
|---|--|----|
| 1 | Report for CNN model for fake review detection | 6 |
| 2 | Report for SVM model for fake review detection | 6 |
| 3 | City frequencies | 10 |
| 4 | Category frequencies | 10 |

List of Figures

| | | |
|---|-----------------------------------|---|
| 1 | SVC truth labels | 6 |
| 2 | Calibrated truth labels | 6 |
| 3 | City distribution | 9 |
| 4 | Category distribution | 9 |

Listings

1 Introduction

1.1 Problem description

We have devised a recommendation system for restaurants based on the Yelp Dataset Challenge.

Our aim is to predict if a certain user will like a certain restaurant, depending on characteristics of the restaurant, user's taste (derived from his previous reviews), opinion of similar users (who gave similar votes to similar restaurants), trustworthiness of the reviews.

In order to do so, we studied three papers, that we used as the basis of our work:

1. Restaurant Recommendation System, Ashish Gandhe [3];
2. Machine Learning and Visualization with Yelp Dataset, Zhiwei Zhang (with her repo) [4];
3. Recommendation for yelp users itself, Wenqi Hou, Gauravi Saha, Manying Tsang [5].

We started from the data cleaning performed by Hou, Saha and Tsang (3), then we applied the SVM model proposed by Zhang (2) to identify fake reviews, and assigned to each review a "truth score" (an indicator of the trustworthiness of the review), so that we could use this score as a weight in the computation of the historical features described by Gandhe (1).

After some further preprocessing we applied three machine learning models to the obtained dataset and got our predictions.

1.2 Datasets

The Yelp dataset we used is available here and described in detail here.

It is composed of five JSON files:

1. *business.json* contains data about businesses (we've narrowed our search to restaurants only, but there were hotels and shops too) including location data, attributes and categories;
2. *review.json* contains full review text, the ids of the user who wrote it and of the restaurant it was about, and the number of stars given as a vote (we used the texts only for the detection of deceptive reviews);
3. *user.json* contains data about users, such as popularity, friends and name;
4. *checkin.json* contains the checkins on a business (we discarded this dataset);
5. *tip.json* contains tips written by a user on a business (we dropped the text and kept only the "compliment count", as a sign of the reliability of a user).

The most significant features are written in or obtained from the review file (2), since the label we want to predict, likes/dislikes, is calculated from the number of stars assigned by a user to a restaurant, and so are some of the historical features presented in Gandhe [3] and discussed in a separate section [2.3].

The review dataset is also the biggest one, with more than 1 million reviews, spanning from 2014 to 2018, and more than 5GB in size.

1.3 Tools used

We summarize and motivate the tools used for this project:

- **Python 3:** the most widely used programming language for machine learning tasks, we used version 3.6 for compatibility with Tensorflow GPU;

- **Jupyter Notebook:** we decided to write our main scripts in a notebook since it's more readable and comments are clearer, even if it gives some problems, with the multiprocessing library for example;
- **Pandas:** we used this library to read and manage the datasets, but it appeared to be inherently inefficient, since it doesn't do anything to overcome the Python GIL, so we tried to parallelize it using Modin and Numba, but also those attempts have failed, so we used the previously mentioned multiprocessing module of the standard library;
- **Scikit-Learn:** it provides ready-to-use implementations of many ML non-neural models;
- **TensorFlow GPU:** it provides the implementations of neural networks.

Thus, most of our code is in the file `/src/main_notebook-multiprocess.ipynb`, the functions that are explicitly run on many processes are in `src/multiproc_utils.py`, the portion of preprocessing customized from [5] is in `/src/recommendation_system_preprocessing.ipynb` and the model for finding deceptive reviews from [4] is in `/Yelp_Sentiment_Analysis/Scripts/fake_reviews.ipynb`.

2 Preprocessing

2.1 Data cleaning

As the first thing we did with our datasets, we just applied the data cleaning by Hou, Saha and Tsang [5] with some adaptation to our purpose, e.g., we want to consider restaurants in all the available cities and not only Las Vegas.

What we do based on their kernel is:

1. Transfer json into Pandas dataframe with proper indexing;
2. Extract data that includes all restaurants;
3. Replace garbage data which includes incorrect states and postal codes;
4. Date transformations and standardization;
5. Create new explanatory features based on initial features;
6. Delete consequently unnecessary columns which could add ambiguity based on new features;
7. Delete duplicate restaurants entries and combine their reviews;
8. Save obtained datasets in pickle format so that they occupy less space in memory.

The main features they work on are **categories**, **attributes** and **hours** in the *business* dataset, renamed *restaurants* since other kind of businesses are not considered.

The feature **categories** contains a list of kinds of businesses, ranging from shops to sport centers, so only the items that contains "restaurant" in that list are kept; then the most frequent and significant categories are extracted and saved in a new feature called **cuisine**, that we used to infer the users' preferences and to compute the historical features [2.3].

The feature **attributes** contains a dictionary of characteristics of the restaurant, such as if it offers alcohol or if it has wifi connection, but these dictionaries are a bit messy, so the authors create a new feature for each attribute and make the values homogeneous, for example **wifi** has values such as "No", "u'no'", "'no'", 'None' and only "No" is kept.

The feature **hours** contains a dictionary of the form {"day_of_the_week": "opening_hour-closing_hour"}, so it's not very usable, thus the authors divide this features in two features per day of the week: one for the opening and one for the closing.

2.2 Fake Review Detection

With our data clean and tidy, we used some models described by Zhang in [4] to detect deceptive reviews and associate a *trustworthiness score* to each of them.

We trained all of their models: a Convolutional Neural Network, fed with sequences obtained by a tokenizer from the Keras preprocessing library, and a Support Vector Machine, fitted on vectors computer by Scikit Learn's TfidfVectorizer.

It turned out that the performances of the Convolutional Neural Network were worse with respect to the Support Vector Machine, as we can see in the tables [1] and [2].

Thus, we used the SVM to compute predictions and establish which reviews are reliable and which are deceptive, in a binary fashion, so we added to the review dataset a feature **bin_truth_score** that holds a 1 for true reviews and a -1 for fake ones.

In order to have a probability distribution instead of a binary result, the SVM is passed as a parameter to a Scikit Learn's CalibratedClassifierCV, whose predictions are the probabilities that a review is authentic, instead of binary results, and we used them to add a feature **real_truth_score** to the review dataset.

Once this was done, we removed the texts of the reviews.

In the figures [1] and [2] we can see the distribution of the two kinds of labels.

| | precision | recall | f1-score | support |
|---------------------|------------------|---------------|-----------------|----------------|
| -1 | 0.18 | 0.18 | 0.18 | 24070 |
| +1 | 0.87 | 0.88 | 0.88 | 158468 |
| micro avg | 0.78 | 0.78 | 0.78 | 182538 |
| macro avg | 0.53 | 0.53 | 0.53 | 182538 |
| weighted avg | 0.78 | 0.78 | 0.78 | 182538 |
| accuracy | 78.3% | | | |

Table 1: Report for CNN model for fake review detection

| | precision | recall | f1-score | support |
|---------------------|------------------|---------------|-----------------|----------------|
| -1 | 0.84 | 0.95 | 0.89 | 161016 |
| +1 | 0.96 | 0.86 | 0.91 | 211105 |
| micro avg | 0.90 | 0.90 | 0.90 | 372121 |
| macro avg | 0.90 | 0.91 | 0.90 | 372121 |
| weighted avg | 0.91 | 0.90 | 0.90 | 372121 |
| accuracy | 89.99% | | | |

Table 2: Report for SVM model for fake review detection

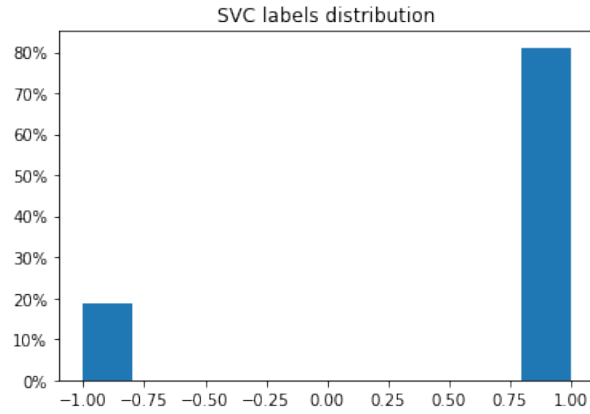


Figure 1: SVC truth labels

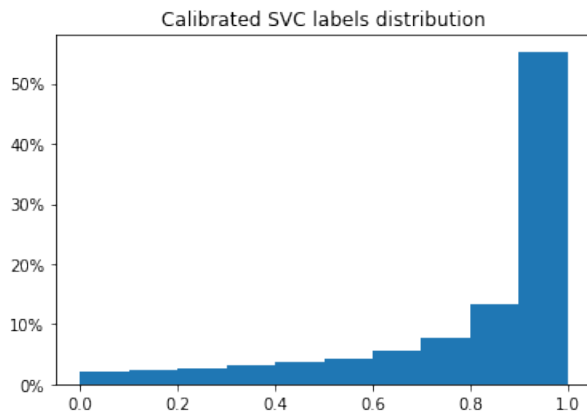


Figure 2: Calibrated truth labels

2.3 Historical features

Based on Gandhe’s paper [3], we have added some *historical features* to our dataset:

1. User-level features:
 - 1.1. average of the ratings given by a certain user,
 - 1.2. number of reviews written by a certain user;
2. Business-level features:
 - 2.1. average of the ratings given to a certain restaurant,
 - 2.2. number of reviews written about a certain restaurant;
3. User-business-level features:
 - 3.1. average rating given by a certain user to each cuisine (feature defined in section Data cleaning), i.e., for each cuisine, the average of all the votes given by a user to all the restaurants that contains that cuisine in their list,
 - 3.2. average of the ratings given by a certain user to the cuisines of a certain restaurant, i.e., the average of the votes given by a user to the cuisines present in the list of a restaurant, as computed in the feature (3.1.).

In order to do this, we had to split the review dataset in three parts:

- *Test set*, from the last day considered in the dataset, to the previous m months;
- *Training set*, from the day before the beginning of the test set, up to n months before;
- *History*, the remaining part of the dataset, used to compute historical features.

We have picked $m = 3$ and $n = 8$, so the test set goes from 9/1/2018 to 11/30/2018, the training set goes from 1/1/2018 to 8/31/2018, the history contains the remaining data, from 10/12/2004 to 12/31/2017.

While the other features were quite straightforward to implement and fast to compute, we had to optimize the code of the feature (3.1.) in order to prevent its computation from taking more days, so we tried two libraries that are supposed to automatically speed up a python program running it on more processors, with a minimal intervention by the programmer. The first, Modin, is specifically meant for Pandas, and the second, Numba, more generically for Python math operations, but none of them actually worked, the former because of some incompatibility with Windows OS, and the latter because it “likes NumPy” (as written in the docs) but not Pandas. So, we decided to manually split the datasets and pass portions of them to a function launched in parallel on multiple processors by a process pool, and this reduced the computation time to just few hours.

We decided to apply our two estimators of the trustworthiness of the reviews (described in Fake Review Detection) to compute two additional version of the just described features, so we have three flavors for each of these features:

1. *standard*: each review is treated in the same way (every review has weight 1),
2. *binary*: the authentic reviews have weight 1 and the deceptive ones have weight 0, accordingly to their `bin.truth.score`, so fake reviews aren’t counted in total or average,
3. *real*: each review is weighted with its `real.truth.score`, both for the total and for the average.

2.4 User-based collaborative approach

Then, we decided to add a new feature based on *collaborative filtering*: the idea is that we have lots of reviews written by lots of users, so we can use the score given by similar user to similar restaurants

to predict what a certain user will think about a restaurant he never tried before. So we applied the following formula:

$$pred(u, r) = a_u + \frac{\sum_{u_i \in U} sim(u, u_i) * (a_{u_i, r} - a_r)}{\sum_{u_i \in U} sim(u, u_i)} \quad (1)$$

where

- a_u is the average of the votes given by the user u to all the restaurants that share some cuisines with r (i.e., the feature 3.2. between user u and restaurant r defined in section Historical features),
- U is the set of all the users in the dataset, u excluded,
- $a_{u_i, r}$ is the same as a_u , but for user u_i , so it is the value of the feature 3.2. for user u_i and restaurant r ,
- a_r is the average rating received by r (i.e., the feature 2.1.),
- $sim(u, u_i)$ is computed as the cosine similarity between two vectors that represent the two users u and u_i , composed by the values for the features 3.1. (one for each cuisine and one for each version, as explained in Historical features).

We have three versions for this feature too, with the same meaning described in the previous section.

2.5 Dimensionality reduction and further preprocessing

rimosso alcuni valori delle features categoriali con troppi valori (perché se no scoppiava la RAM), riempito i null con media o moda, convertito features categoriali in numeriche.

At this point we had all the features we needed, so, that was the moment of preparing the data so that they were ready to be fed to the models. The next steps we performed were the following:

1. We joined all the interesting datasets in only one, to have all the needed features together and one review per row (the *checkin* dataset was discarded because not useful);
2. We added the label **likes** that holds a 1 if that user gave 4 or 5 stars in that review or 0 if he/she gave 1, 2 or 3 stars;
3. We removed unnecessary columns such as **review date**, **restaurant name**, **restaurant address**, **user name**;
4. We applied some dimensionality reduction (more details on this later);
5. We filled missing values with the mode of the feature for the categorical values and with the mean for the numerical ones;
6. We converted categorical features into numerical features, readable by the models:
 - for the features **OutdoorSeating**, **BusinessAcceptsCreditCards**, **RestaurantsDelivery**, **RestaurantsReservations**, **WiFi**, **Alcohol**, **city** we used Pandas' `get_dummies()` function, that applies one hot encoding on the input features,
 - for the features **Monday_Open**, **Tuesday_Open**, **Wednesday_Open**, **Thursday_Open**, **Friday_Open**, **Saturday_Open**, **Sunday_Open**, **Monday_Close**, **Tuesday_Close**, **Wednesday_Close**, **Thursday_Close**, **Friday_Close**, **Saturday_Close**, **Sunday_Close**, **postal_code** we applied Scikit learn's `OrdinalEncoder.fit_transform()`, that encodes categorical features as an integer array: it could bias the models but it is less memory consuming with respect to one hot encoding, so we preferred this method for features that we consider less determinant and that have many possible values.

Now let's describe what we have done for dimensionality reduction. We had two features that seem to be quite important and that could assume hundreds of values: **city** and **categories**. First of all, we plotted the distribution of the values of these features, as shown in the figures [3] and [4]: it appears that there are few very frequent values (such as 'Las Vegas', 'Phoenix', 'Toronto' for the cities

and ‘Food’, ‘Nightlife’, ‘Bars’ for the categories) and lots of values that are pretty rare, thus we decided to cut the less frequent values before proceeding with one hot encoding, so we choose a threshold $\theta = 100$ and replaced all the values that occur less than θ times with ‘other’.

Let’s note that, since after one hot encoding we have one column for each category, the previously created **cuisine** (see Data cleaning section) has become useless, so we removed it.

In the tables [3] and [4] are listed the ten most frequent and ten least frequent cities/categories.

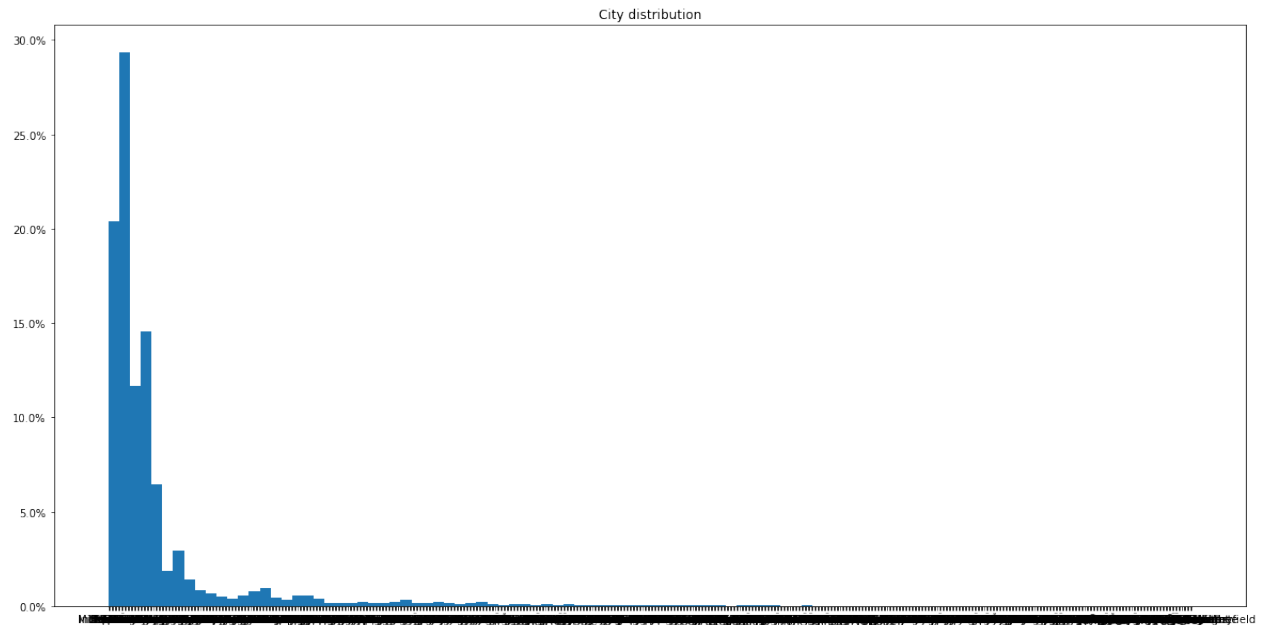


Figure 3: City distribution

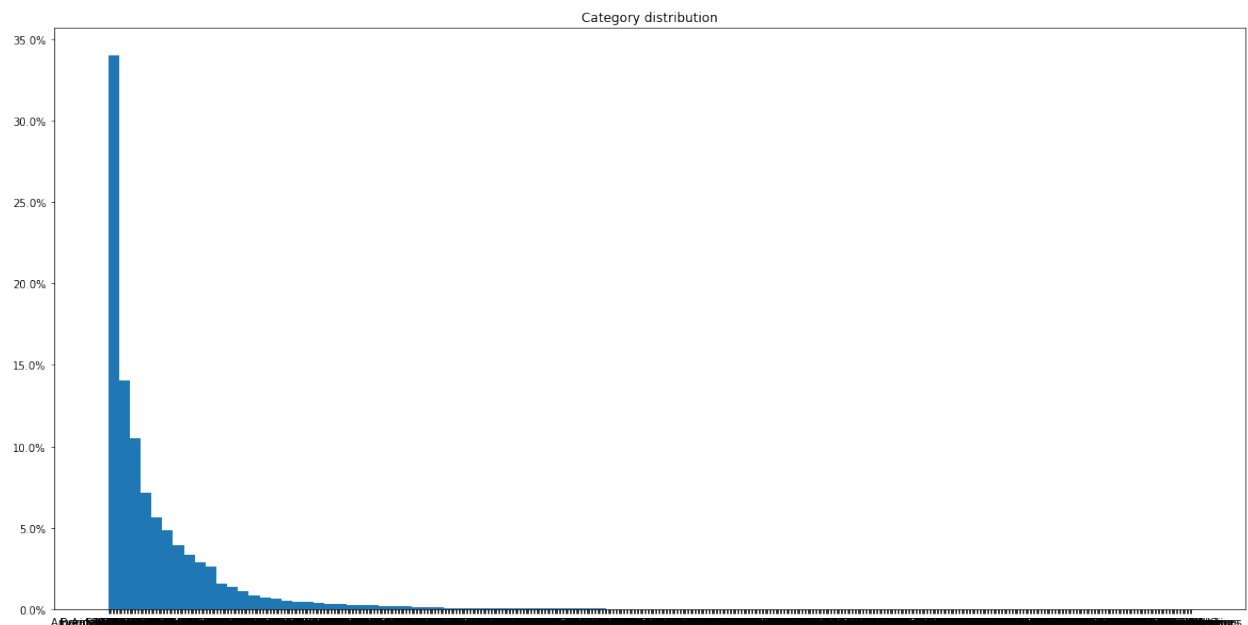


Figure 4: Category distribution

| city | # occurrences |
|---------------------|----------------------|
| Las Vegas | 208437 |
| Phoenix | 71126 |
| Toronto | 57047 |
| Charlotte | 40190 |
| Scottsdale | 36579 |
| Pittsburgh | 26891 |
| Henderson | 21518 |
| Montreal | 18531 |
| Tempe | 18354 |
| Mesa | 17235 |
| ... | ... |
| Paw Creek | 1 |
| Tottenham | 1 |
| springdale | 1 |
| Coteau-du-Lac | 1 |
| Huntingdon | 1 |
| Fabreville | 1 |
| De Winton | 1 |
| Napierville | 1 |
| Laval, Ste Dorothee | 1 |
| Les Coteaux | 1 |

Table 3: City frequencies

| category | # occurrences |
|------------------------|----------------------|
| Food | 192841 |
| Nightlife | 170259 |
| Bars | 165959 |
| American (Traditional) | 123975 |
| Breakfast & Brunch | 120310 |
| American (New) | 117437 |
| Sandwiches | 82264 |
| Mexican | 73726 |
| Burgers | 71598 |
| Pizza | 67538 |
| ... | ... |
| Beer Hall | 1 |
| Banks & Credit Unions | 1 |
| University Housing | 1 |
| Pet Groomers | 1 |
| Gardeners | 1 |
| Home Health Care | 1 |
| Campgrounds | 1 |
| Holiday Decorations | 1 |
| Roofing | 1 |
| Senegalese | 1 |

Table 4: Category frequencies

3 Models

Qui descriviamo solo i modelli, motivandone la scelta.

3.1 Linear Support Vector Machine

Per confronto col paper

3.2 Random Fores

Per avere un esempio di ensemble

3.3 CNN???

Per avere un esempio di reti neurali

```
#Questo e' un commento  
for i in lst:  
    print("ciao")
```

4 Experiments and evaluation

Qui descriviamo tutti i test che abbiamo fatto per ogni modello, con parametri e risultati.

4.1 Linear Support Vector Machine

- addestramento su 1/2 training set, 5000 iterazioni, nessun filtro città o categorie → best training score 0.732, best test score 0.701
- addestramento su 1/3 training set, 10000 iterazioni, nessun filtro città o categorie → best training score 0.725, best test score 0.736
- addestramento su 1/3 training set, 10000 iterazioni, filtro città più di 100 occorrenze, filtro categorie più di 200 occorrenze → best training score 0.722, best test score 0.706
- addestramento su training set completo, 50000 iterazioni, filtro città e categorie più di 100 occorrenze, nuova feature basata su score di utenti simili → best training score 0.743, best test score 0.737

mai raggiunta convergenza

4.2 Random Forests

grid search con 900 parametri

4.3 CNN???

bla bla

5 Conclusions

Siamo stati bravi dacci 30

6 References

1. Yelp dataset on Kaggle <https://www.kaggle.com/yelp-dataset/yelp-dataset>;
2. Yelp Dataset presentation <https://www.yelp.com/dataset>;
3. Restaurant Recommendation System, Ashish Gandhe
<https://www.semanticscholar.org/paper/Restaurant-Recommendation-System-Gandhe/093cecc3e53f2ba4c0c466ad3d8294ba64962050>;
4. Machine Learning and Visualization with Yelp Dataset, Zhiwei Zhang
https://medium.com/@zhiwei_zhang/final-blog-642fb9c7e781
(with her repo https://github.com/zzhang83/Yelp_Sentiment_Analysis);
5. Recommendation for yelp users itself, Wenqi Hou, Gauravi Saha, Manying Tsang
<https://www.kaggle.com/wenqihou828/recommendation-for-yelp-users-itself>;