

Modelling Rating Prediction Algorithms on Google Local Reviews: A Comparative Study

Abstract—Google Reviews have become an integral part of our daily lives. We decide places to visit, restaurants to dine, and services to utilize based on these reviews. Furthermore, the general population and their sentiments about a business have a huge impact on the performance of the business. In this project, we utilize the Google Local Reviews data set to predict the potential rating when a user reviews a business. First, we describe our pre-processing and feature engineering strategies for extracting relevant information from the data set. We then proceed to define the predictive modelling problem and subsequently train several models along with baselines. Finally, we analyze the results obtained and conclude by providing a summary of the models' performances. We also discuss the relevant literature.

Index Terms—Google Reviews, Regression, Machine Learning, Deep Learning, Data Modelling, Personalized Recommendation.

I. INTRODUCTION

The Google Local Reviews Data set¹ consists of three files that contain information specific to users, businesses, and reviews, respectively, along with metadata. To derive meaningful information from the data and to define a predictive problem, we first clean, transform, prune, and split the data. Then, we perform exploratory data analysis and data visualization to find interesting trends, patterns, and seasonalities in the local reviews data. We perform feature engineering on the pruned data set, thereby creating useful features for optimizing the regression problem at hand. Here, we intend to predict the rating a user would give to a business based on the user and the business/place metadata. Additionally, we look at a scenario where we do not have these metadata attributes, i.e., we have only the (userID, placeID) pairs. We approach both these tasks by proposing several baseline (intuitive) models and the design of a few machine and deep learning models, along with the evaluation metrics. Our aim is

to create a model that fits the data and generates precise predictions².

II. LITERATURE REVIEW

The data set we use for our analysis and experiments is the raw version of the Google Local Reviews data. We transform the data in such a way that we can prepare regression-like predictive problems and optimize the representation of the data set through encoding techniques, process and clean the raw data and also bring about notions of data completeness, sparseness and cardinality in the process of analyzing and pre-processing the data. Several regression and classification based problems exist in literature for such data sets. For example

Similar data sets studied in the past include x,y,z.

Some of the state-of-the-art models found for predictive tasks involving such data sets include, xx,yy,zz....

We in turn propose a problem statement on the data set in such a way that we can employ some of these techniques and we devise a couple of our own data pre-processing methods to prepare the data such that it optimizes regression performance and makes the data much more intuitive and meaningful. Subsequent sections describe how we achieve the same.

III. EXPLORATORY DATA ANALYSIS & INSIGHTS

In this section we describe the exploratory data analysis carried out on the Google Local Reviews data set (hereinafter referred to as 'data set').

A. Data Pre-processing

We perform the following pre-processing on the data set:

1. We merge the 'user' and 'places' data tables from the data set based on the criterion of 'gPlusUserId' and 'gPlusPlaceId' matching and we obtain a resulting 'merged' table. We drop all the columns that have extremely low completeness or extremely high randomness

¹https://cseweb.ucsd.edu/~jmcauley/datasets.html#google_local

²Our code can be found at <https://github.com/AwsManas/Assignment-2>

and uniqueness as they would skew our data and could result in sub-optimal models. The resulting dataframe is shown below. The dataframe has 8649011 rows and 11 columns.

| merged.head() | | | | | | | | | | |
|---------------|-------------------|-----------------------|--------|--------------------------------------|-------------------|-----------------------|----------------|-------------------|------------------------|------------|
| | userName | gPlusUserId | rating | reviewText | categories | gPlusPlaceId | unixReviewTime | name | price | gps closed |
| 0 | an lam | 10000010817154263736 | 3.0 | Chất lượng tạm ổn | [Già Trĩ - Cafe] | 108103314380004200232 | 1.372867e+09 | Cà Phê Thăng Long | [10.852044, 106.65971] | False |
| 1 | hoang long nguyen | 101658842775082396018 | 5.0 | Good coffee, nice and peaceful place | [Già Trĩ - Cafe] | 108103314380004200232 | 1.354888e+09 | Cà Phê Thăng Long | [10.852044, 106.65971] | False |
| 2 | Hong Le | 107574984242995460712 | 2.0 | Chợ heo | [Già Trĩ - Cafe] | 108103314380004200232 | 1.352015e+09 | Cà Phê Thăng Long | [10.852044, 106.65971] | False |
| 3 | HALL TURGUT | 100000013500285534661 | 5.0 | We si temiz duzenli.. | [Turkish Cuisine] | 102194128241608748649 | 1.342871e+09 | Selale Restaurant | [37.8037, 29.2209] | False |
| 4 | Akudorot Yazlim | 105271324704942360981 | 5.0 | None | [Turkish Cuisine] | 102194128241608748649 | 1.373148e+09 | Selale Restaurant | [37.8037, 29.2209] | False |

Fig. 1. Merged Dataframe

2. We prune the data set by dropping columns with excessive NaN values, and perform NaN imputation on relevant columns like 'gps'. We also expand the gps column to two columns ['lat', 'lon'] indicating latitude and longitude (geolocation) of the place. We thus have a cleaned data set which we now call the 'pruned' version of the data with 1299688 rows \times 13 columns.

3. We utilize this pruned data set and geolocations to create a mapping of the various types of ratings that we see across all businesses/places in the world. The pruned data set and its properties are given below.

B. EDA & Outlier detection

We perform the following on the pruned data set to identify meaningful patterns and extrapolate from the data.

1. Compute the basic stats of every column in the data set and understand their relevance for modelling a predictive task. We check the following components:

- Completeness of the entire data set and that of every column in the data set.
- Missingness of the entire data set and that of every column in the data set.
- Distinctness of the entire data set and that of every column in the data set.

2. We compute the following stats in the data set:

- Mean value
- Standard deviation
- Min, Max Values
- Percentiles
- Entropy

We infer that columns with high entropy (or degree of randomness) can be removed and the max and min values gave us a sense of how we can normalize the data. Data scaling and pruning was performed further with the information we obtained through EDA.

```
In [20]: final_df.isna().sum()
Out[20]:
userName      0
gPlusUserId    0
rating         0
reviewText     2649295
categories     135813
gPlusPlaceId   0
unixReviewTime 42698
name           0
price          6782945
gps            27431
closed         0
dtype: int64

In [21]: dfSize = pruned_df.size

In [22]: ## Ratios

In [23]: null_sums = final_df.isna().sum().sum()

In [24]: #Completeness of entire dataframe
null_sums / dfSize * 100

Out[24]: 57.04435336670149

In [25]: # Missingness of entire dataframe
(dfSize - null_sums) * 100 / dfSize

Out[25]: 42.9556466332985
```

Fig. 2. Some Computed Statistics

```
In [31]: from math import e

def pandas_entropy(column, base=None):
    vc = pd.Series(column).value_counts(normalize=True, sort=False)
    base = e if base is None else base
    return -(vc * np.log(vc) / np.log(base)).sum()

In [34]: pandas_entropy(pruned_df['rating'])
Out[34]: 1.4216977063054153

In [43]: pandas_entropy(pruned_df['unixReviewTime'])
Out[43]: 13.92883155136558

In [33]: pandas_entropy(pruned_df['reviewTextLength'])
Out[33]: 6.423664854916784
```

Fig. 3. Some Computed Statistics

3. Next we perform some outlier detection in our pruned data set to remove skewness-causing entries in our data. This will greatly help us in regularization of our models.

Examples to depict some of the outlier detection techniques we employed are:

- a. Outlier detection based on z-score thresholding: Here we use our 'reviewTextLength' feature (discussed in a later section) and we compute the z-score for every entry of that feature. We retain only the rows where $-2 < z_score < 2$ for the given feature. This way we remove outliers i.e. for ex:- Users who give only 1-10 character reviews or give 1000-2000 character reviews as many of these are junk reviews and will bias our model extensively.
- b. Cardinality based pruning: We drop columns with low & high cardinalities based on some thresholds we defined for each.
- c. Percentile based outlier detection: We remove all the reviews/data points that have their review-TextLength on the extremely higher side of the percentile distribution i.e. 99-100 percentile range.
- d. Outlier filtering based on timestamp: We noticed

that there were some extremely old and irrelevant/out of trend reviews such as a review from 1990 and a review from 2001 etc while the most recent review for the specific place was circa 2014. It makes sense to keep only the most recent reviews and not reviews that were given 2 decades ago.

C. Data Visualizations & Inference

Once we completed our data pruning, outlier detection, scaling and normalization, we visualized various stats and columns in the data so as to derive more meaningful patterns that we can exploit in our models. This information is subsequently used in data set localization and for feature engineering purposes.

Given below are the screenshots of several visualizations that we performed on the data set and we provide all of our insights subsequently.

```
In [16]: plt.hist(pruned_df['rating'],bins = 5)
plt.show()
```

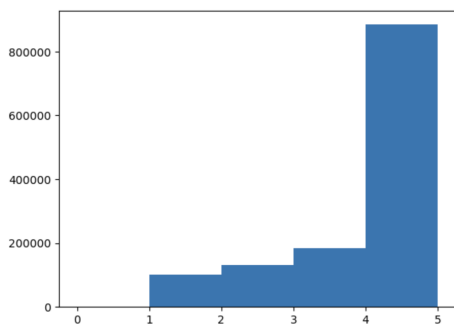


Fig. 4. Rating distribution cross-population

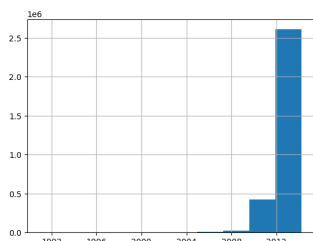


Fig. 5. No. of Ratings available year-wise

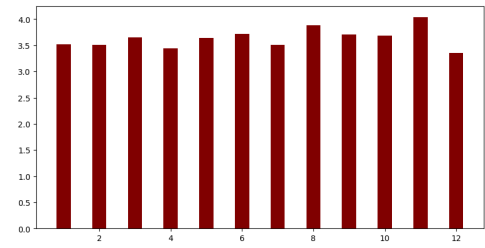


Fig. 6. Month Wise Rating Distribution in U.S.

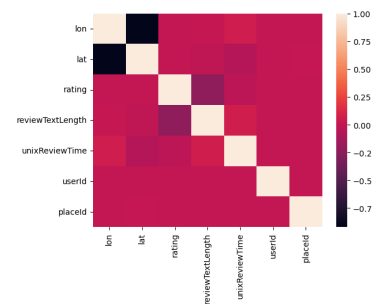


Fig. 7. Correlation Map of features for all data

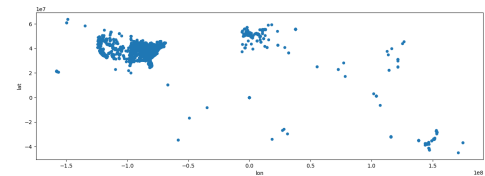


Fig. 8. World Map of Locations present in Places data

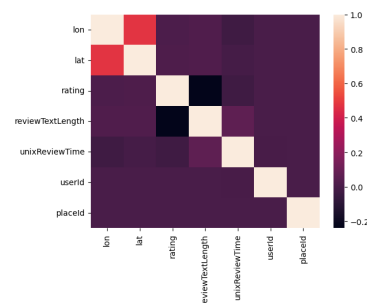


Fig. 9. Correlation Map of features of ratings data across the US

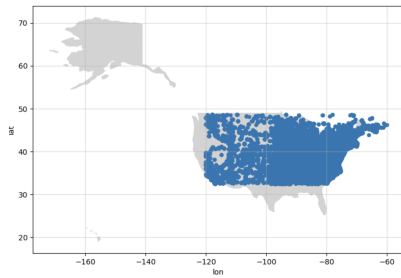


Fig. 10. Localized reviews distribution across U.S.

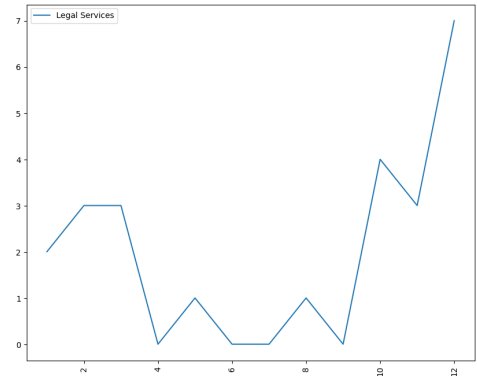


Fig. 14. Monthly Trends in Rating of Legal Services

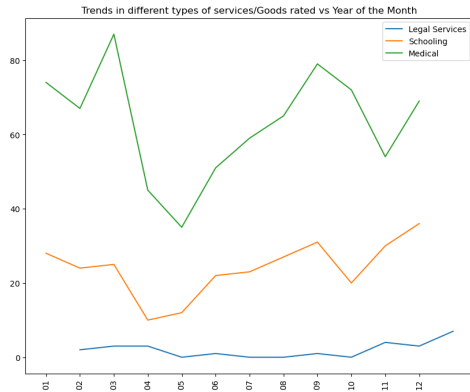


Fig. 11. Monthly Trends in Rating of Various Services

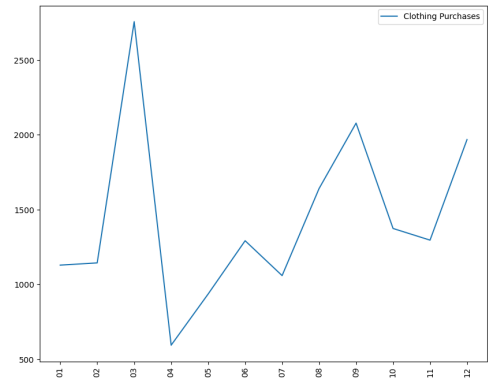


Fig. 15. Monthly Trends in Rating of Cloth Purchases



Fig. 12. US state-wise average rating across all services

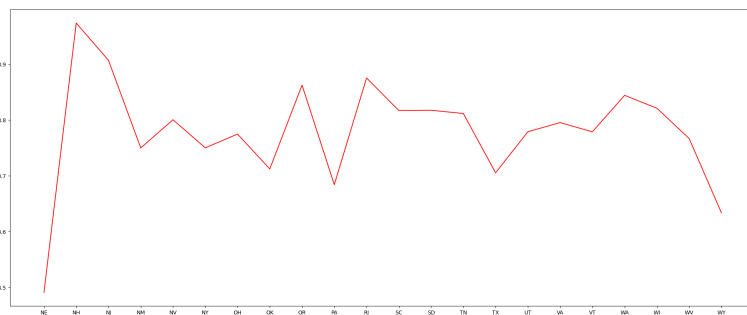


Fig. 13. US state-wise average rating across all services

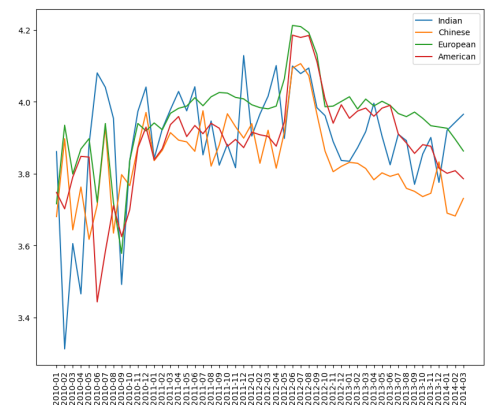


Fig. 17. Temporal rating trends in 4 Types of Restaurants

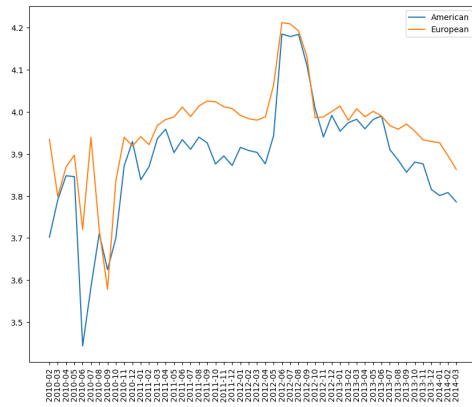


Fig. 16. Temporal rating trends in American & European Restaurants



Fig. 18. User rating trends for a subset of users

1) Analysis and Insights from the Visualizations: :

Apart from some obvious preliminary intuitions, we've collated some deeper insights and correlations that we observed in the data in this section as a prelude to feature engineering.

- Majority of the ratings >62.5% tend to be on the higher side of 4-5 stars. This indicates that users tend to be more biased towards rating higher than lower.
- We see from the world map distribution that there is a dense set (8 million) of ratings in the United States alone, which we can further localize into states to obtain state-wise business/place ratings.
- The categories column has a strong correlation on the type of people that rate a business and the influence of categories are dynamic over the years.
- We see from temporal visualizations that the ratings tend to vary drastically based on the time of the year that a particular business is rated. This means that the 'year-month' data point would be a strong indicator of user ratings.
- Categories can be clubbed. For Ex:- a specific place can have multiple categories associated with it, but

on an n-dimensional vector space, we can see that these categories lie close enough that they can be grouped to a single category thereby dimensionality reduction can be employed.

- The location of a specific business has a huge impact on the business rating, performance (i.e. whether the business would close or not) and also on user sentiments. The variation can be seen in multiple graphs.
- It is also evident from the visualizations that though the sentiments/ratings of 2 individuals for the same place/business in a large populas may vary drastically based on their tastes. The general rating trend for two countries (or large populations) as a whole are collective and are similar in distribution. The ratings are consistent over time.
- Demand, Utilization and rating of certain services or businesses such as 'Gift Shops', 'Hospitals', 'Schools', 'Legal Services' etc. vary greatly based on the time of the year, and a trend in the data is noticed across countries as well.

Considering these insights, we perform some feature engineering on the data set next.

D. Feature Engineering

The following describe our feature engineering steps:

- First we reduce the scope of our data set to account for only interactions (i.e. user-business ratings) that occurred within the United States. This gave us about 3 million data points to work with. The way we extracted this information was on the basis of the latitude and longitude of the GPS data.
- We wrote python scripts to find the Geospatial distance in miles from the <latitude, longitude> of our business/place to the center of each state in the US, and used the nearest neighbour approach to zero-in on which state in the US a particular business lies in. Thus we extracted a new piece of relevant information as an indicator of the rating given to a place.
- Since it was evident from our temporal correlation analysis that the column 'unixReviewTime' is highly relevant, we extract the 'month-year' information from it and create 'month' and 'year' features.
- Our EDA showed that the 'categories' column in the dataframe had extremely high cardinality but since the 'categories' column is too important of a feature to drop we resorted to a customized

category grouping technique using 'word' similarities (i.e. of all the categories that a particular place/business belongs to, which category best describes the place). We posit that there exists 10 such mutually exclusive categories that we can optimally 'bin' our places into. which are ('Associations/Organizations', 'Entertainment', 'Legal Services', 'Medical', 'Public', 'Restaurant', 'School', 'Shops and Stores', 'Venues', 'Others') and we bin all the data points into each of these categories. Thus a new high-correlation feature is created.

- Additionally, Our category wise user rating prediction analysis in the visualizations showed a strong dependency of average user rating for a category with the places that belong to those categories. Therefore, we systematically generated the "User Average Ratings" for every business grouped "Category-Wise" and included that as another data point in the parameter vector.
- Finally we perform a one-hot-encoding of the columns 'year', 'month', 'state', 'final_category' to get the final data set.

IV. PREDICTIVE MODELLING & TASKS IDENTIFICATION

Before we define our predictive model, here we attach a glimpse of what our final data set looks like after all the steps in Section 2.

Based on this, Some of the possible predictive tasks for this data set include:

1. Given user and place metadata (i.e. all columns in the data set except rating), predict what rating a user would give to a business.
2. Given the location, and historical user-ratings of a place, predict whether a business would fail/succeed.
3. Predicting the top-k most similar users for a particular user given the way they rated different places in different categories.

In this report we discuss the implementation, experiments and analysis of Task-1. For modelling Task-1 we first concretize our problem statement, and suggested line of thinking. We narrow down the various approaches we can use and algorithms that fit the task, and discuss how we should evaluate our model.

Problem Statement: For a user with rating characteristics 'x' and for a place with metadata such as 'location (state)', 'month', 'year' etc. develop a model that predicts what rating the user would give to the place.

We propose regression-based techniques to develop solutions to this problem as we need to make real-valued predictions in the range 1-5 (star-ratings). Multi-variate regression is intuitively the ideal way to approach this problem. So Initially we start by building a naive baseline and then suggest intelligent baselines. We later on model the regression problem and perform a comparative analysis of the performance of several models for this task. The performance is evaluated by choosing suitable evaluation metrics. We also attempt to avoid over-fitting by exploring various techniques of regularization.

V. MODEL SELECTION & EXPERIMENTS

The Experimental setup we use are the following: 1. Data is split into 3 subsets 'train', 'validation', 'test' in the ratio 60:20:20. 2. Regularization is done wherever applicable ensuring to optimize the validation MSE. 3. Experiment epochs are set such that they converge when the difference in inter-trial MSE's are < 0.0001 . 4. We progress from simpler to more personalized and complicated models and evaluate them across with the same baselines.

The models that we chose, with parameters and a couple of evaluation metrics we considered for our experiments are as follows:

A. Baselines

1) *Always predict mean:* Our first naive baseline is predicting the mean rating of all businesses for all users satisfying any condition such as falling within a specific state/category. The different variants of this baseline are to Always predict mean:

- Per state: Find avg of all ratings per state and return the average rating for every state as the predicted rating.
- Per category: Do the same but based on category
- Per unit of time: Do the same but based on month and year of the rating.

2) *Always predict popular mean:* An alternative baseline can be to Always predict the mean rating of the most popular place/business satisfying a specific condition. The variants of this baseline are:

- Per state: predict the average rating of the most popular place in that state.
- Per category: Do the same but for the most popular place of that category.
- Per unit of time: Do the same but for the most popular item within that month-year time frame.

Table I below summarizes the baseline performances.

TABLE I
BASELINE PERFORMANCE ON TEST SET

| Baseline Model | Mean Absolute Error | Mean Squared Error |
|----------------------------------|---------------------|--------------------|
| Always Predict Mean | | |
| Per State | 0 | 1.674 |
| Per Category | 0 | 1.673 |
| Per unit of time | 0 | 1.662 |
| Always Predict Popularity | | |
| Per State | 0 | 2.335 |
| Per Category | 0 | 1.802 |
| Per unit of time | 0 | 2.688 |

B. Models

- Linear Regression
- Ridge Regression
- Fix regularization alpha in Ridge regression
- Random Forest Regression
- XGBoost Regression
- Bayesian Personalized Ranking (i.e., without user/item metadata)
- BPR with latent factors
- Tensorflow DNN

TABLE II
REGRESSION MODEL PERFORMANCE

| Model | Test MAE | Test MSE |
|-----------------------------------|----------|----------|
| Linear Reg. | 0 | 0.4172 |
| Ridge Reg. | 0 | 0.4172 |
| Ridge with reg. ³ | 0 | 0.4172 |
| Random Forest | 0 | 0.4245 |
| XGBoost | 0 | 0.4173 |
| BPR w/ LF <u,i> only ⁴ | 1.5905 | 1.0235 |
| BPR w/ LF <metadata> ⁵ | 0 | 0 |
| TF-DNN ⁶ | 0.3654 | 0.4166 |

C. Evaluation Metrics

- MAE
- RMSE / MSE
- R-square / Adjusted R-square

The table below depicts the performance of each of these models on the data set. We consider the MSE to be the evaluation metric.

VI. RESULTS & CONCLUSION