

# The Yelp Dataset: Topic Modelling and Text Classification

Lorenzo Lecce n. matr. 830881

Andrea Lucini Paioni n. matr. 826578

## ABSTRACT

This project takes inspiration from a part of the Yelp dataset, a dataset which includes 6,990,280 reviews from 150,346 businesses, sourced from the online platform Yelp, famous for sharing opinions on local businesses. The study employed topic modeling and text classification techniques to unveil the primary themes and content within the reviews, categorizing them into distinct groups. Initial steps involved text preprocessing steps, including normalization, stopwords removal, tokenization, lemmatization... Latent Dirichlet Allocation (LDA) was chosen for topic modeling, while two text representations (TF-IDF and Doc2Vec) were evaluated for the text classification task. The findings revealed the identification of 10 food-related topics through topic modeling. In terms of text classification, few model were developed, with the best one achieving an accuracy of 87%, a recall score of 96% and an F1 score of 90%, in addressing a classification problem with the aim of predicting the score of the reviews based on the text content of the reviews.

## CONTENTS

<b>ABSTRACT</b>	<b>1</b>
<b>CONTENTS</b>	<b>1</b>
<b>INTRODUCTION</b>	<b>1</b>
<b>OBJECTIVES</b>	<b>2</b>
<b>DATA PREPARATION</b>	<b>2</b>
<b>TEXT PREPROCESSING</b>	<b>2</b>
<b>EXPLORATORY ANALYSIS</b>	<b>3</b>
<b>TOPIC MODELING</b>	<b>4</b>

<b>TEXT CLASSIFICATION</b>	<b>5</b>
Classification with <i>countvectorizer()</i>	5
Classification with TFIDF	5
<b>RESULTS AND CONCLUSIONS</b>	<b>5</b>
<b>TABLE INDEX</b>	<b>6</b>
<b>FIGURES INDEX</b>	<b>6</b>
<b>REFERENCES</b>	<b>6</b>

## INTRODUCTION

Yelp serves as an online platform where individuals share their perspectives on local businesses, offering a space to review various establishments like restaurants, local businesses, and hotels. Users can explore and evaluate businesses across different categories. Moreover, businesses have the opportunity to create and manage their profiles on Yelp, showcasing information about their offerings, along with photos and contact details. The platform facilitates user reviews, which can be filtered based on criteria such as rating, date, and location.

Yelp extends tools to businesses, aiding them in managing their online presence and monitoring the performance of their listings. Widely embraced, Yelp is considered a valuable resource for both consumers and businesses. The complex and huge dataset provided by the platform is composed by different datasets:

- "yelp\_academic\_dataset\_review": a dataset containing 6,990,280 reviews and information on the businesses on the platform.

- *"yelp\_academic\_dataset\_business"*: a dataset containing info about the businesses on the platform.
- *"yelp\_academic\_dataset\_user"*: a dataset containing info about the users who write reviews on the platform.
- *"yelp\_academic\_dataset\_tip"*: a dataset containing info about some particular reviews about the businesses on the platform.
- *"yelp\_academic\_dataset\_checkin"*: a dataset containing info about the time of the reviews about the businesses on the platform.

In this project, the focus is on employing topic modeling to extract the primary themes and content from reviews. Additionally, classification techniques are applied to categorize texts of the reviews. With these tasks in mind, we thought that only the *"yelp\_academic\_dataset\_review"* and the *"yelp\_academic\_dataset\_business"* dataset were needed for our purposes.

## **OBJECTIVES**

The research questions for this project are:

- Evaluate the performances of different text representations considering the classification task.
- perform topic modelling techniques, to find some of the most discussed topics in the Yelp reviews.
- predicting the review stars considering the text of the reviews with classification.

## **DATA PREPARATION**

The datasets considered (*"yelp\_academic\_dataset\_review"* and *"yelp\_academic\_dataset\_business"*) were huge, so we had to consider a sample, and we reduced also the number of variables, because some of them weren't useful for our purposes.

In fact, the reviews dataset had 6990820 instances (with 9 variables), and the businesses dataset had 150345 observations (with 14 variables). We decided to consider only 7

variables from the businesses dataset, and then we merged the two datasets, bringing the total to 15 variables:

- *review\_id*: to univocally identify the id;
- *user\_id*: to univocally identify the user who wrote the review;
- *business\_id*: to univocally identify the business;
- *review\_stars*: the stars assigned to the review by the user;
- *useful*: how many users found the review useful;
- *funny*: how many users found the review funny;
- *cool*: how many users found the review cool;
- *text*: the full text of the review;
- *date*: date and time of the review;
- *name*: the name of the business;
- *city*: the city of the business;
- *stars*: the mean stars of the business reviews, in categories between 1 and 5 stars (every 0.5 stars);
- *review\_count*: how many reviews the business have;
- *attributes*: multiple attributes associated with the business, in a list format;
- *categories*: a list of tags associated with the business;

However, we still had the problem of computational power needed to process a dataset made of six million instances; in order to solve that, we decided to sample the dataset, keeping only the reviews of businesses with the category *"Italian Restaurant"*, bringing the total down to 439358 instances.

## **TEXT PREPROCESSING**

The text processing phase was crucial in our project, because it ensured the following analysis to be more accurate and significant. In particular, all the following phases were performed on the variable *text*, which contains all the text from the customers reviews.

In the preprocessing phase of our dataset, several essential steps were undertaken to ensure its readiness for analysis. Initially, all text within the reviews was converted to lowercase to maintain uniformity. Following that, numerical values were removed, ensuring that the dataset focused solely on textual content.

Tabs, empty lines (e.g., `\n`), and links or URLs embedded in the reviews were systematically eliminated to enhance the cleanliness of the text. Furthermore, white spaces, emojis, and repeated characters were systematically stripped away to refine the data. The removal of punctuation, tokenization, and subsequent elimination of stopwords contributed to the streamlining of the text for subsequent analyses.

Finally, lemmatization was performed, aiming to reduce words to their base or root form, thereby facilitating a more cohesive and standardized representation of the textual information. These meticulous preprocessing steps collectively played a crucial role in preparing the dataset for robust and meaningful text mining endeavors.

## EXPLORATORY ANALYSIS

In our dataset analysis, we conducted three distinct examinations to capture comprehensive insights (all the following analysis were conducted on a copy of the business dataset only, filtered to consider only Italian restaurants). The initial analysis involved generating a bar plot to illustrate the distribution of the stars from the ratings.

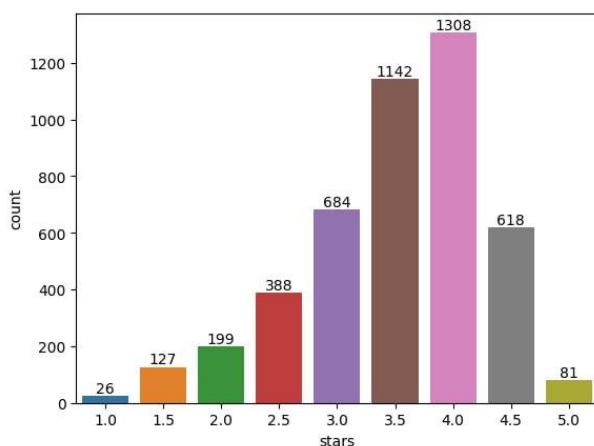


Figure 1: barplot variable "stars"

We noticed that the ratings of 3.5 and 4 were the most prevalent categories, meaning there was a general concentration of positive sentiments within the reviews.

The second analysis was focused on the review count across different categories, with a plot providing information about the number of reviews for each business activity considered in our dataset.

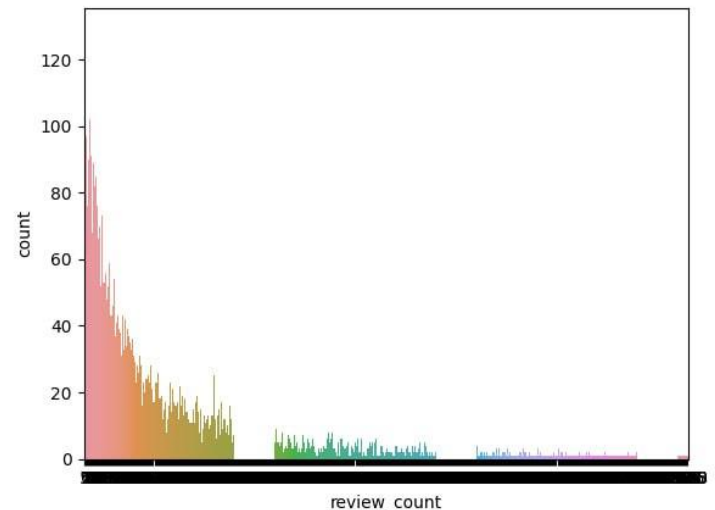


Figure 2: barplot variable "review\_count"

Lastly, we developed a boxplot to present the distribution of the variable "stars" (the same of Figure 1) within the dataset. Together, these analyses helped us with interesting insights about significant variables from the dataset, giving useful info on both the sentiment distribution and review frequency patterns.

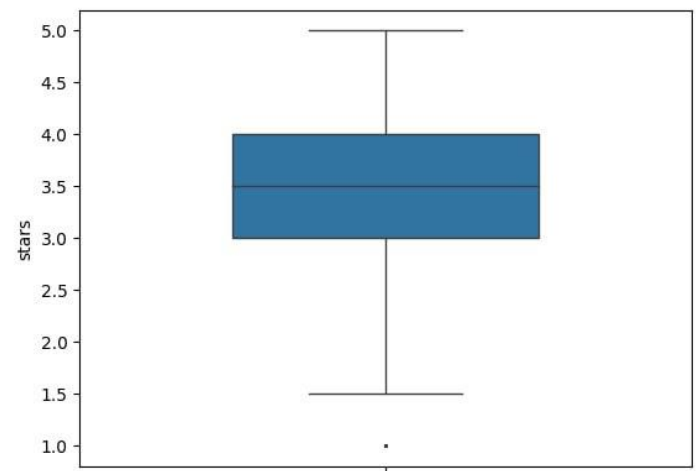


Figure 3: boxplot variable "stars"

<i>var</i>	<i>mean</i>	<i>min</i>	<i>max</i>	<i>std</i>	25%	50%	75%
<i>Stars</i>	3.51	1.0	5.0	0.78	3.0	3.5	4.0
<i>Review count</i>	92.39	5	4250	156.77	18	44	108

We confirmed the information we had obtained with the plots: the distribution of *"stars"* is asymmetric towards higher values (all the values between 25% and 75% are between 3.0 and 4.0), with a mean of 3.51; and for the *"review\_count"* variable we have the opposite situation, with a mean of 92.39 and a distribution towards lower values (75% of the business have less than 108 reviews).

To ensure the correctness of the task we decided to develop, we decided to further sample the dataset: some task on both topic modeling and text classification required a lot of time and computational power.

The goal of the topic modeling phase was to identify a series of topics included in the text from the reviews of the businesses.



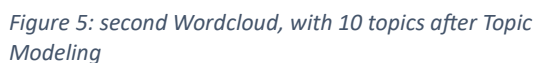
As we can see, some of the most common words used are expected, and related in particular with Italian food: “pizza”, “pasta”, and some comments probably referred to the quality of Italian food, such as “good”, “great”, “well”, “nice”, “delicious”...

It uses a Bayesian inference system to estimate the most common topics in the text collection, with also a weight associated with each topic to show their importance.

An important aspect of the LDA model is the fact that we need to choose an appropriate number of topics: we decided to consider 10 topics. We obtained a model with 10 different topics.

Given the fact that the categories of restaurants were very similar, it's understandable that the topics were very similar with each other: in general, "pizza" and "good food" were the most prominent topics found, with "great time", "come

The 10 topics described can be seen in the Wordcloud reported below:



The text classification task is a mix of supervised methods aimed at classifying a characteristic of the text data; it can be both a multi-class classification and a binary outcome one.

For the classification purposes, we considered the variable “*stars*”, that indicates the number of stars, between 1 and 5, given by the customer to each business in association with the text review. We created another variable, called “*good\_bad*”, with a binary outcome: “1” if the review was considered good (only 4 or 5 stars), “0” if the review was considered bad (either 1, 2 or 3 stars).

We divided the dataset in train set (90% of the data) and test set (10% of the data), for the purpose of classification.

These below are the metrics obtained as a result of the classification methods performed on data represented with `countvectorizer()` representation:

Table 2: results classification with countvectorizer() representation for Decision Tree, Support Vector Machine and Random Forest classifier

These below are the metrics obtained as a result of the classification methods performed on data represented with *TFIDF()* representation:

Table 3: results classification with TFIDF() representation for Decision Tree, Support Vector Machine and Random Forest classifier

Our project focuses on two task, topic modeling and text classification, performed on the Yelp dataset.

For the topic modeling task, we expected to see some words related to Italian food as the most common topics, and our expectations were confirmed: “pizza” and “pasta” were among the most common words, and in particular “pasta” was among the most important words in each of the 10 topics obtained; other positive words were among the most common in the topics obtained, such as “good”, “great”, “come back”, “great service” and “good time”.

For the purposes of text classification, we considered three different classification methods (a Decision Tree classifier, a Support Vector Machine classifier and a Random Forest classifier) and two text representations (a *TFIDF* representation, and a *countvectorizer* representation): we expected better results for the *TFIDF* representation, as it focuses both on the frequency of words in the corpus and it also provides the importance of words. Our expectations were confirmed.

Among the different classification methods, the Support Vector Machine (with *TFIDF*, only method able to get all scores above 0.90 on the metrics) classifier achieved the best results metric-wise, while the Decision Tree obtained the worst results among the methods chosen.

- <https://www.yelp.com/dataset/>
- <https://scikit-learn.org/stable/>
- <https://towardsdatascience.com/latent-dirichlet-allocation-lda-9d1cd064ffa2>
- <https://www.jmlr.org/papers/volume3/blei03a/blei03a.pdf>

## **TABLE INDEX**

<i>Table 1: statistics about the variables "stars" and "review_count"</i>	4
<i>Table 2: results classification with countvectorizer() representation for Decision Tree, Support Vector Machine and Random Forest classifier</i>	5
<i>Table 3: results classification with TFIDF() representation for Decision Tree, Support Vector Machine and Random Forest classifier</i>	5

## **FIGURES INDEX**

<i>Figure 1: barplot variable "stars"</i>	3
<i>Figure 2: barplot variable "review_count"</i>	3
<i>Figure 3: boxplot variable "stars"</i>	3
<i>Figure 4: first general Wordcloud</i>	4
<i>Figure 5: second Wordcloud, with 10 topics after Topic Modeling</i>	5

## **REFERENCES**

- Text Mining and Search notes