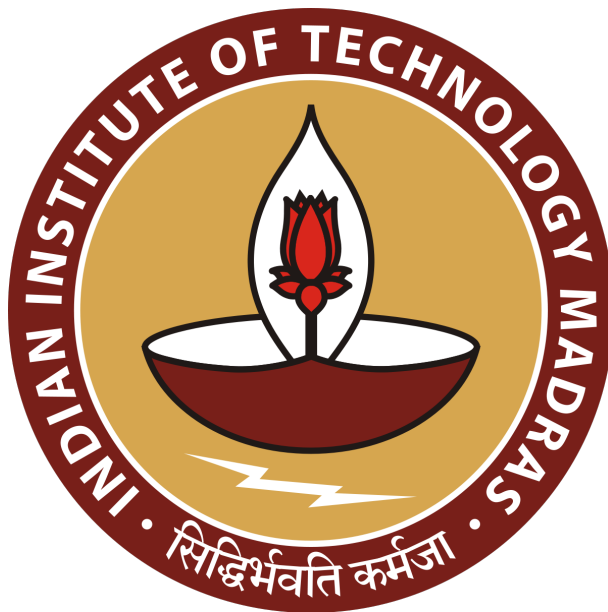# CS4830 - BIG DATA LAB
## FINAL PROJECT

---

# Classification on the Yelp dataset using NLP

---

AKASH ANAND - ED18B002

NIHAL S MANVI - AE18B031

SHASHANK H S - BE18B006

INDIAN INSTITUTE OF TECHNOLOGY MADRAS

JAN - MAY 2022

# Contents

# List of Figures

# 1 Introduction

In this project, we are provided with a real-world dataset upon which we perform the required analysis using tools and services learnt during the course on the Google Cloud Platform. The dataset that we are provided with is a subset of data from the Yelp website which contains many user-submitted reviews of various businesses.

# 2 Problem Statement

The dataset provided to us is over 3 GB large. Our aim is to classify a review into a given star-rating category by performing sentiment analysis using NLP on the review text.

Our approach was to pre-process the provided data and train several machine learning models on it to come up with the predictions. We use Kafka streaming model for testing the model and making predictions in the subscriber on data coming in from the publisher. The required steps that have been performed in this project are as follows:

1. Exploratory Data Analysis

2. Preprocessing and Feature Engineering

3. Training of different models

4. Accuracy/F1-score comparisons

5. Kafka streaming

# 3 Exploratory Data Analysis and Data Preprocessing

Default image version for cluster is Debian 10, Hadoop 3.2, Spark 3.1. Due to Anaconda's unavailability with these specifications, we changed the cluster image version to Ubuntu 18.04 LTS, Hadoop 2.10, Spark 2.4.

A schema of the dataset is shown in the figure below.

```
The original yelp dataset has 6600446 rows.
root
 |-- id1: string (nullable = true)
 |-- ufc1: float (nullable = true)
 |-- date: string (nullable = true)
 |-- ufc2: float (nullable = true)
 |-- id2: string (nullable = true)
 |-- stars: float (nullable = true)
 |-- text: string (nullable = true)
 |-- ufc3: float (nullable = true)
 |-- id3: string (nullable = true)
```

Figure 1: Schema of Yelp Data

Some more details about the various features in the provided csv follow. Three input labels are given as id1, id2, id3 which represent three unique identifiers:

1. User ID: Unique ID provided by Yelp to a particular user

2. Business ID: Unique ID provided by Yelp to a particular business

3. Review ID: Unique ID provided by Yelp to a particular review

With the provided csv file, it is unclear which id column corresponds to which exact id.

Three more input labels ufc1, ufc2, ufc3 are provided which are numerical ratings for:

1. Useful - The number of people who found the review useful

2. Funny - The number of people who found the review funny

3. Cool - The number of people who found the review cool

Again, it is unclear with the provided csv file as to which column corresponds to which type of rating. The above uncertainties in column name does not affect training in any capacity.

Additional input labels which are clear as per the given csv file are:

1. Date - Date on which the review was uploaded

2. Text - Text feedback provided by the user as part of the review

The output label that has been provided to us is the star rating: the rating of a particular business given by an user.

A summary of the dataset provided to us is:

Original yelp dataset

| | id1 | ufc1 | date | ufc2 | id2 | stars | text | ufc3 | id3 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | --30_8lhuyMHbSOcNWd6DQ | 0.0 | 2018-04-30 16:56:24 | 0.0 | DBJXMz1Rlr91eFUd0gKjrA | 4.0 | My child has been attending since he was 4. H... | 1.0 | LH2QmlXtq8CLji65INs76Q |
| 1 | --7PUidqRWpRSpXebiyxTg | 0.0 | 2019-11-02 23:30:02 | 0.0 | 0O4RqD91ZHdCPw4rUzBIYA | 1.0 | Lunchtime while visiting a family member at th... | 0.0 | IzuszJGVrtdxda-sKXwh0Q |
| 2 | --8lbOsAAxjKRoYsBFL-PA | 0.0 | 2015-06-16 02:10:33 | 0.0 | UYYEUH8qldJg4C4XqVX3Gg | 2.0 | Eh.. not that great don't waste your money on ... | 6.0 | gdcRlubKDmslUYFPHUp1Cg |
| 3 | --OS_l7dnABrXvRCCuWOGQ | 0.0 | 2019-03-14 15:52:59 | 0.0 | gOChygeYyXUhL7D0Krw_hw | 5.0 | We recently had our 93 Cadillac repaired at Le... | 0.0 | DzWToz-VRSOGB6Cje57NoA |
| 4 | --eBbs3HpZYlym5pEw8Qdw | 0.0 | 2019-03-29 22:31:53 | 0.0 | NjlC7n0vlRFqx9EToHvNUg | 1.0 | Very poor experience. We had 5 queen size room... | 0.0 | Um1ayg0vDsEzWRDVwqM8zA |

Figure 2: Yelp Data Before Processing

Original yelp dataset summary
```
+-------+-----------------+------------------+-----------------+------------------+
|summary|             ufc1|              ufc2|             ufc3|             stars|
+-------+-----------------+------------------+-----------------+------------------+
|  count|          5330044|           5191554|          5056834|           5191535|
|   mean|0.6932374274442319|0.29718461948002467|1.082546526146558|3.8076801947786154|
| stddev|60.225810914976144| 1.8224370589136485|6.418358538911953|1.452101366531364|
|    min|        -75.54586|              -1.0|             -1.0|               0.0|
|    max|          75614.0|            2019.0|           2021.0|               7.0|
+-------+-----------------+------------------+-----------------+------------------+
```

Figure 3: Summary of Yelp Data Before Processing

The data in the above figure came with many rows having many of the features having NaN values. As a result, such rows were naturally filtered from the dataset. A summary of the obtained data after processing is shown below:

```
The processed yelp dataset has 5056831 rows.
Yelp dataset summary with ufc1, ufc2, ufc3 and stars null value rows removed
```

| | id1 | ufc1 | date | ufc2 | id2 | stars | text | ufc3 | id3 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | ---kPU91CF4Lq2-WIRu9Lw | 0.0 | 2020-01-29 18:39:02 | 0.0 | 8NnKwxC71uLNWs00efgD7w | 5.0 | Unfortunately the site for Frankie's is incorr... | 0.0 | YhbCO1DVINYkmVv8DCAIxw |
| 1 | ---kPU91CF4Lq2-WIRu9Lw | 0.0 | 2020-02-20 00:50:49 | 0.0 | XawsDBeNXIT_rRivcgmhyA | 5.0 | Never would have thought a little hidden place... | 1.0 | 5jIO2REcgB6GKFeSsc-OXw |
| 2 | ---kPU91CF4Lq2-WIRu9Lw | 0.0 | 2020-12-14 01:17:21 | 0.0 | s7f2L3EESkKf-kNDSchpow | 4.0 | I love this place. Nice place for the girls to... | 1.0 | goqGoC76zemDagYvRa8GlA |
| 3 | ---kPU91CF4Lq2-WIRu9Lw | 0.0 | 2021-04-12 19:46:10 | 0.0 | HwGGdjfpl7_ndf9d9W-6fw | 1.0 | Be careful before ordering the coleslaw, I fou... | 1.0 | YUVbBNr_dSJNP2pwDk1xyg |
| 4 | ---kPU91CF4Lq2-WIRu9Lw | 0.0 | 2021-10-17 02:02:58 | 0.0 | gANpst_byMcYH6c6nqRcRg | 5.0 | What a fantastic casual seafood or BBQ local s... | 0.0 | jtPb5gfrvYixrC0axWjqlA |

Figure 4: Yelp Data Post Processing

```
Processed yelp dataset summary
+-------+------------------+-----------------+-----------------+------------------+
|summary|              ufc1|             ufc2|             ufc3|             stars|
+-------+------------------+-----------------+-----------------+------------------+
|  count|           5056831|          5056831|          5056831|           5056831|
|   mean|0.4602334940598173| 0.278760947320565|1.082545981860933|3.8256057993632773|
| stddev| 2.060862884241416|1.5443312626601509|6.418360280637526|1.4447704407498172|
|    min|              -1.0|             -1.0|             -1.0|               1.0|
|    max|             404.0|            378.0|           2021.0|               5.0|
+-------+------------------+-----------------+-----------------+------------------+
```

Figure 5: Summary of Yelp Data Post Processing

Some overarching insights we gain from the above data follow:

1. Star values range from 0 to 7 in integer increments before processing, and 0 to 5 after processing.

2. There is no particular limited range as such for useful, funny, cool variables because the value completely depends on other users on the website and their response to the reviews

3

## 3.1 Graphical Representations



Figure 6: Bar Plot Distribution of Stars Data

From the above plot, we can see that majority of the reviews have 5-star ratings while very few have 2-star ratings.

Figure 7: Bar Plot Distribution of ufc1 Data



Figure 8: Bar Plot Distribution of ufc2 Data

The plots of ufc1 and ufc2 show that the data is highly skewed towards the lower end of the X axis.



Figure 9: Bar Plot Distribution of ufc3 Data



Figure 10: Correlation Matrix of the Features

From this correlation matrix, we can see that the given numerical input variables ufc1, ufc2,

ufc3 have very low correlation with 'stars' variables. Thus, from this we can sense that there will be high dependency on the 'text' data variable for classification modelling

# 4 Pre-processing and Feature Engineering

The following series of operations were performed on the 'text' column to obtain a tf-idf vector representation of the field.

1. Document assembler - Transforms the given data into the required format for the Annotator in NLP Spark to work on

2. Tokenizer - It breaks the raw text into small chunks. Tokenization breaks the raw text into words called tokens. These tokens help in understanding the context or developing the model for the NLP. The tokenization helps in interpreting the meaning of the text by analyzing the sequence of words. For example, the text "It is raining" can be tokenized into 'It', 'is', 'raining'

3. Normalizer - Removes all dirty characters from text following a regex pattern and transforms words based on a provided dictionary. It brings all text to lowercase and removes punctuation.

4. Stemmer - Works by cutting off the end of the word, taking into account a list of common suffixes that can be found in an inflected word. This indiscriminate cutting can be successful in some occasions, but not always.

5. Finisher - This is used to close the annotation and transformation done for the NLP spark module's usage

6. Stopword Remover - This removes all the stopwords like 'is', 'a', 'an', 'not', 'are', etc in the given text. This cleaner dataset helps in feature engineering

After the pre-processing the text as described above, feature engineering was done on the text using TF-IDF and Word2Vec. Since Word2Vec took too long, TF-IDF was chosen for feature engineering.

# 5 Training Models

After the pre-processing and feature engineering, we use the modified dataset to train logistic regression, Naive Bayes, decision tree models. Since the dataset is huge and the TF-IDF vectors obtained after preprocessing are large, choosing a complex model (such as SVM or Random Forest) would potentially result in very high training time. Therefore, we used simple models such as Naive Bayes Classifier, Logistic Regression and Decision Tree Classifier. The training is initially done on a small part of the dataset and it is observed as to which model is performing best on the subset of the data. Later, the model giving the highest accuracy is selected for training on the complete dataset.

For training, we used a Dataproc Cluster with one master node and two worker nodes (n1-standard-2 worker and n1-standard-4 master).

The results of training are given below:

## 5.1 Logistic Regression



Figure 11: Results of Logistic Regression

## 5.2 Naive Bayes



Figure 12: Results of Naive Bayes

## 5.3 Decision Tree



Figure 13: Results of Decision Tree

## 5.4 Discussion of Results

A summary of the results from the above models is shown below:

| Accuracy | | |
|---|---|---|
| Model | Training Set | Test Set |
| Logistic Regression | 0.861 | 0.590 |
| Naive Bayes | 0.057 | 0.112 |
| Decision Tree | 0.391 | 0.386 |

From the above results, it is evident the Logistic Regression model performs best on the train and test dataset. We use this model to proceed training the final dataset.

## 5.5 Final Model

Proof of training the final model is shown below:

Figure 14: Training Proof 1



Figure 15: Training Proof 2

Figure 16: Training Set Result



Figure 17: Test Set Result

From the above two figures, we can see that the training set F1 score is 0.7165 and 0.6616 for the test set. This model took over 8 hours to train and evaluate. We could not try different hyperparameter values to improve accuracy before the report deadline.

# 6  Kafka Streaming

For the Kafka streaming task, we have created the publisher and subscriber. Code for these tasks can be found attached. The publisher and subscriber run simultaneously. The subscriber takes the data that is provided by the publisher. Upon receiving the data, the subscriber runs the model and predicts the output which can be seen in the log of the subscriber job.

## 6.1  Publisher

A kafka VM was created and the publisher code was run on the virtual machine SSH window. It is shown as follows:

Figure 18: Publisher

We can see from the figure above that the output in the window corresponds to each data point and the output shows up in the JSON format.

## 6.2 Subscriber

We establish a subscriber on a separate Dataproc cluster simultaneously. As the subscriber code runs, the specifics of the job can be found below:

The outputs are seen in batches. So as we update the publisher (or when it does while going over the dataset), each review is taken in one batch and the accuracy is calculated cumulatively over the tested reviews.

To test the real-time nature of this environment, we input our own review in the required 'json' format in the SSH window of the publisher. This input is then taken by the subscriber to run the model over and predict the stars rating.

The results for some different batches are shown below:

Figure 19: Subscriber 1



Figure 20: Subscriber 2

Figure 21: Subscriber 3



Figure 22: Subscriber 4

Figure 23: Subscriber 5

While testing the model on streaming data, the model was usually found to have an accuracy in the range 60-70% (with some outliers) and an F1 score of around 0.50-0.55.

# 7 Challenges

Some challenges we faced during the course of this project are enumerated below:

1. The provided CSV file contained commas (within the quotation marks) which were not to be inferred as delimiters. Furthermore, new lines were inserted in between quotes causing the data to be read incorrectly into the dataframe.

2. Few problems were faced in utilizing the spark NLP module in the codes. This eventually led a lot of time spent in identifying and utilizing the correct jar files and correct versions of pyspark, spark-nlp and kafka to be used. Hence, we had to use structured streaming and we ran the sub.py file on a new window of the kafka VM.

# 8 Conclusion

1. The given dataset of Yelp is highly skewed. The numerical variables have very less correlation with the 'stars' (output) variable.

2. The logistic regression model works the best for the given dataset.

3. The skewness of the model plays a major role and because of the large computation time, the accuracy is low for the given dataset. Since the 5 stars in the rating are not exactly orthogonal, we feel that this should be treated as a regression problem.

4. Kafka streaming was executed properly as shown above and we were able to obtain accuracies on the test data in real time.

5. The training set accuracy is 71.65% and the test set accuracy is 66.16%.