# Yelp Restaurant Review Analysis with Big Data Technologies

**Nushrat Jahan Ria**
nria1@lsu.edu
**Louisiana State University, Baton Rouge, LA**

## Project Description

- It focuses on the analysis of Yelp's Business Dataset for predicting business success by taking into consideration customer reviews, star ratings and review counts, together with location data applying big data technologies.

- Data is processed in a distributed way using PySpark. PySpark offers efficient processing and transformation at scale through the DataFrame API that makes working with structured data easy.

- The train and evaluate some machine learning models, such as Logistic Regression, Stochastic Gradient Descent, Gradient Boosting, and Random Forest on business performance prediction. Thus, the integration of Dask within this scope will be affirmative to efficient parallel computing and memory optimization for tasks that need scaling beyond local memory constraints, complementing PySpark's processing.

- Cassandra is employed for the distributed storage system, whereas Kafka is used to allow real-time data streaming, thus laying the bedrock for real-time analytics of continuously updating Yelp reviews.
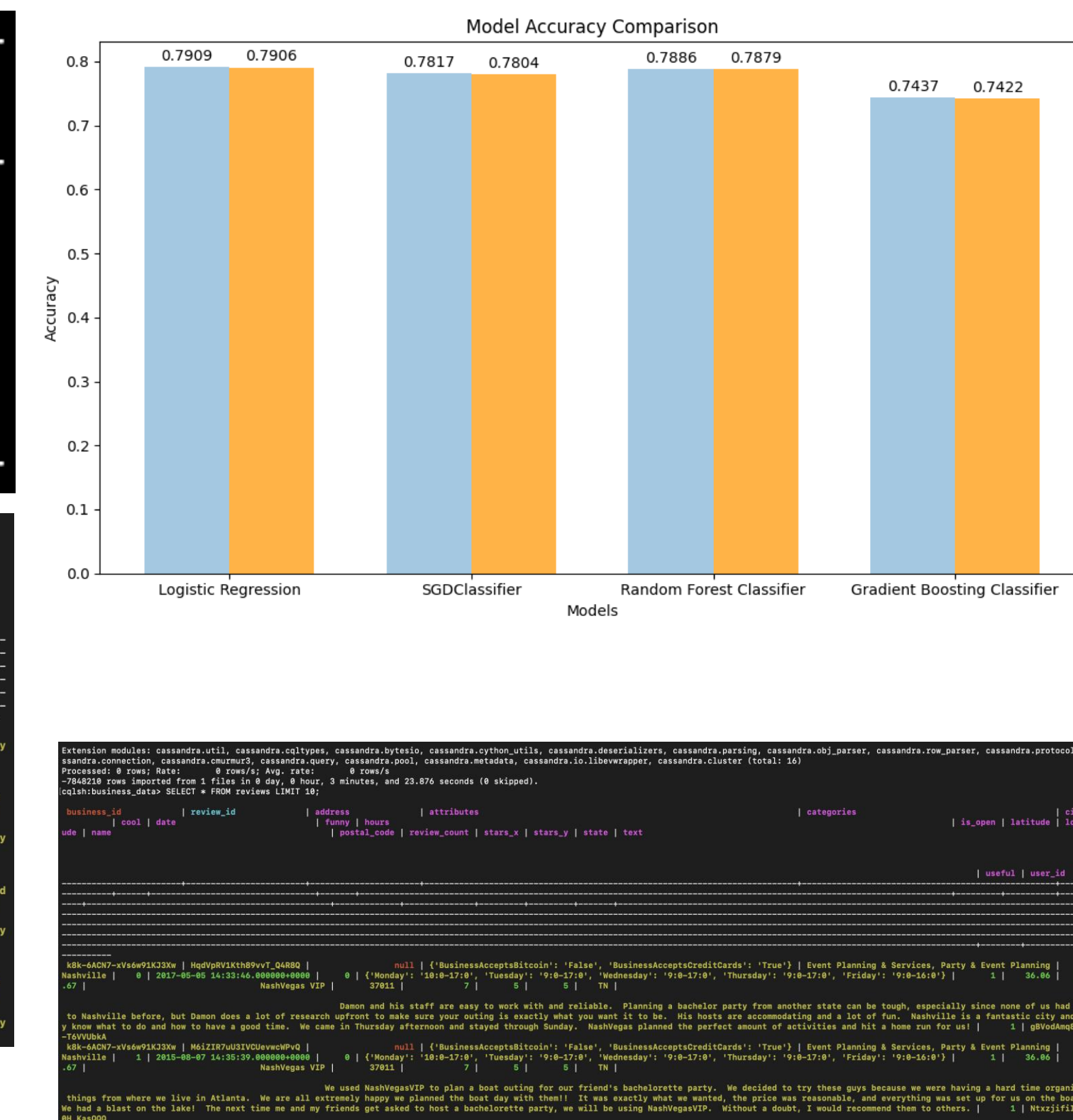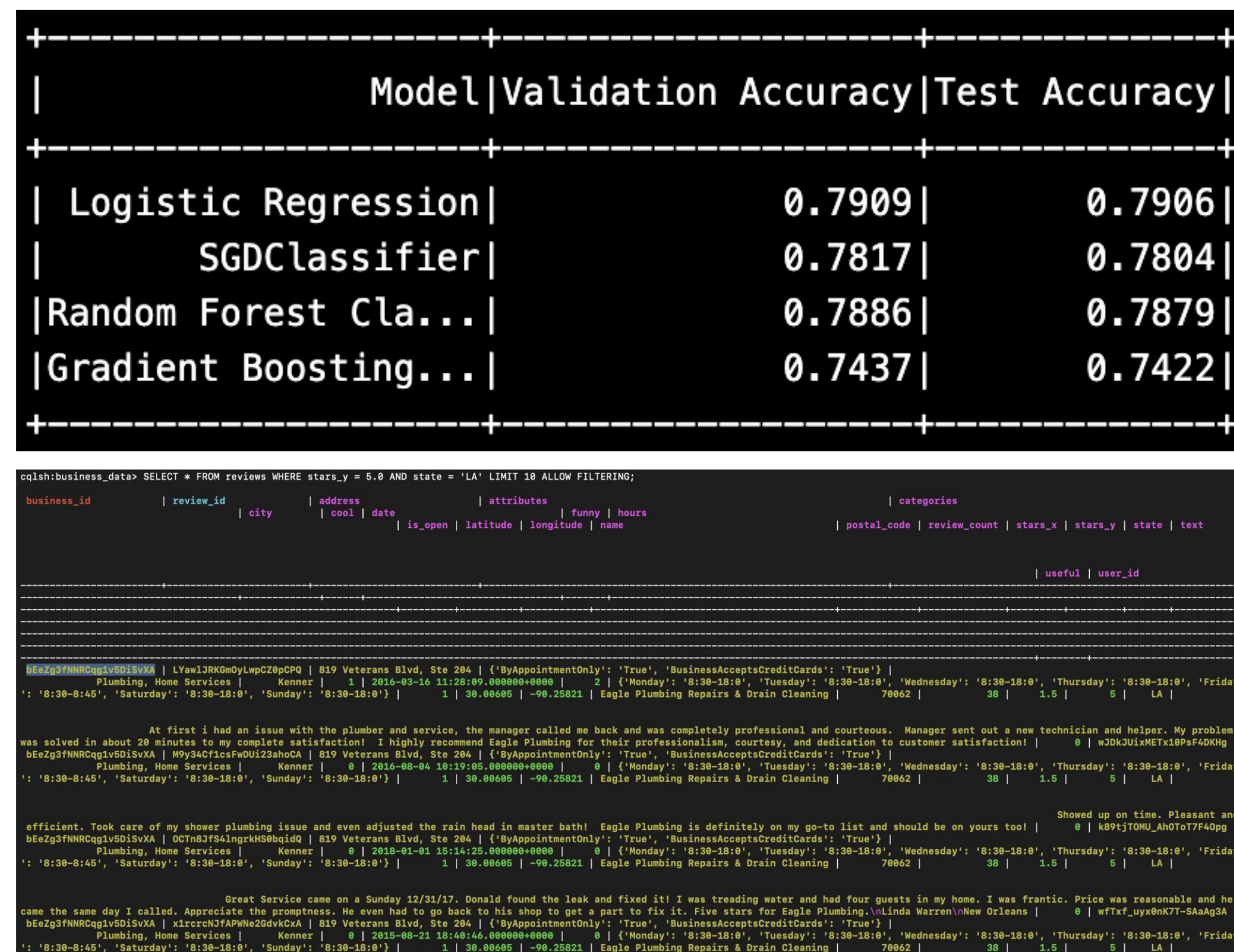
## Development Details

- The project therefore begins by ingesting and cleaning the data; hence, the Yelp business dataset shall be loaded into PySpark DataFrames through missing value handling, renaming columns, and preparing the data for analysis.

- PySpark is used to conduct EDA that explains, in particular, the main features regarding star rating distribution and the counts of reviews to bring out trends and insight into the performance at businesses.

- It designs a machine learning pipeline that trains the models, such as Logistic Regression, SGD, Gradient Boosting, and Random Forest, on the data for the prediction of business success based on customer feedback.

- The performance of the models is estimated by accuracy, precision, recall, and F1-score evaluation metrics, where the Random Forest and Gradient Boosting ensemble methods perform better.

- Big data frameworks are integrated into this system, including PySpark, Dask, Cassandra, and Kafka, where data could be scalably processed, efficiently stored, or even streamed in real time for future development.

## Big Data Frameworks

- **PySpark:** This is designed for large-scale distributed data processing, fast in-memory computation, and parallel processing across clusters; hence, it is indispensable in efficiently processing large-scale datasets.

- **Dask:** This is for parallel computation when tasks do not fit into memory. Integrated with PySpark, it allows for scalable data analysis through efficient handling of heavy computation.

- **Cassandra:** This is a distributed NoSQL database intended to provide scalable data storage, featuring high availability and fault tolerance. The architecture is fitting for this project since it very efficiently stores and retrieves large volumes of structured business data.

- **Kafka:** This becomes useful for real-time data streaming. In fact, Kafka allows for the continuous ingestion and processing of streams of data, thus laying the bedrock for any future real-time Yelp review analysis.

- **Cassandra Driver:** Cassandra is a Python library that interfaces with Cassandra and allows for efficient access and modification of data within the distributed storage system. This especially plays an important role when working with big data sets, such as those provided by Yelp.

## Main Results & Screenshots





## Datasets

- This is the business dataset, a sub-dataset of the Yelp Academic Dataset; for every type of business, this dataset shows restaurants, stores, service providers, and more.

- It includes important features such as star ratings given overall, review counts, business categories, and location details such as city, state, country.

- The data is both structured-some numerical fields comprise rating and review count-while other data is more categorically oriented, such as the type of business and its location. It is ideal for exploratory data analysis and machine learning.

- The dataset is in JSON format. Later, I convert it into CSV format and load it into PySpark for distributed processing, ensuring efficient handling and transformation of large volumes of business data.

- It contains information for several thousand enterprises and is thus a good basis for training machine learning models and doing predictive analytics on business success.

  - https://www.yelp.com/dataset

## Conclusion & Future Work

- The project has effectively leveraged big data frameworks like PySpark, Dask, Cassandra, and Kafka in processing and analyzing large Yelp data. In addition, it has identified Random Forest and Gradient Boosting as the most appropriate predictions of business success.

- Future improvements include hyperparameter tuning, with the exploration of deep learning methods for further improvement in the accuracy of predictions and enabling more complex interactions between features.

- The extension of the system to real-time data processing using Apache Kafka and its integration with sentiment analysis using NLP will provide dynamic insights and a deep understanding of customer feedback.