

Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

Tools

Modeling

conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl



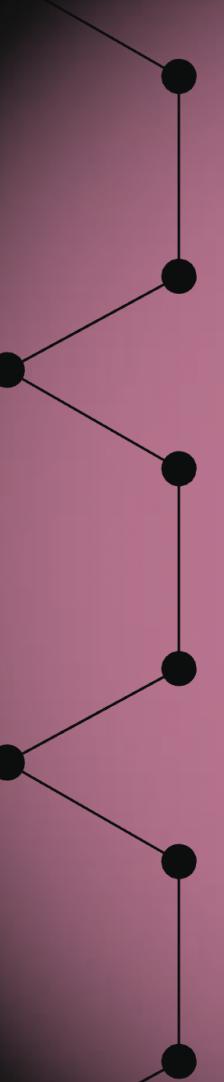


## Description and Objectives

Description

Objectives

The Data &  
Data source/s



What is our project about?



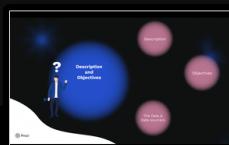
## Description and Objectives

Description

Objectives

The Data &  
Data source/s

- **Get the data**
- **EDA**
  - Data cleaning.
  - Visualization to know the relations between features.
- **Modeling**
  - Data preparation.
  - Apply different classification algorithms.
  - Compare between the results.



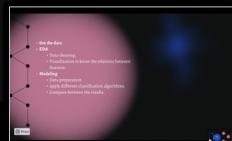


## Description and Objectives

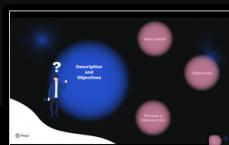
Description

Objectives

The Data &  
Data source/s



- Source: Kaggle.
- 4 Features .
- 100,000 rows, and each row represents a review.





## Description and Objectives

Description

Objectives

The Data &  
Data source/s



Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

Tools

Modeling

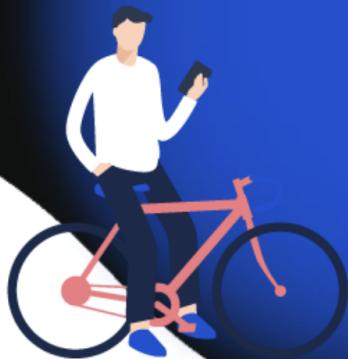
conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl



## Data Preparation



feature  
selection

data cleaning

Feature  
reduction

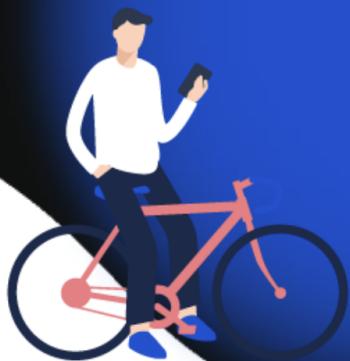
The final  
dataset

## Feature Selection

Dropping all columns except text and label



## Data Preparation



feature  
selection

data cleaning

Feature  
reduction

The final  
dataset



## Data Cleaning

- Tokenize.
- Remove punctuations and numbers.
- Remove stop words .
- Convert to lower case.
- Use stemmer
- Create Bag of Words.
- Compute TF.



## Data Preparation



feature selection

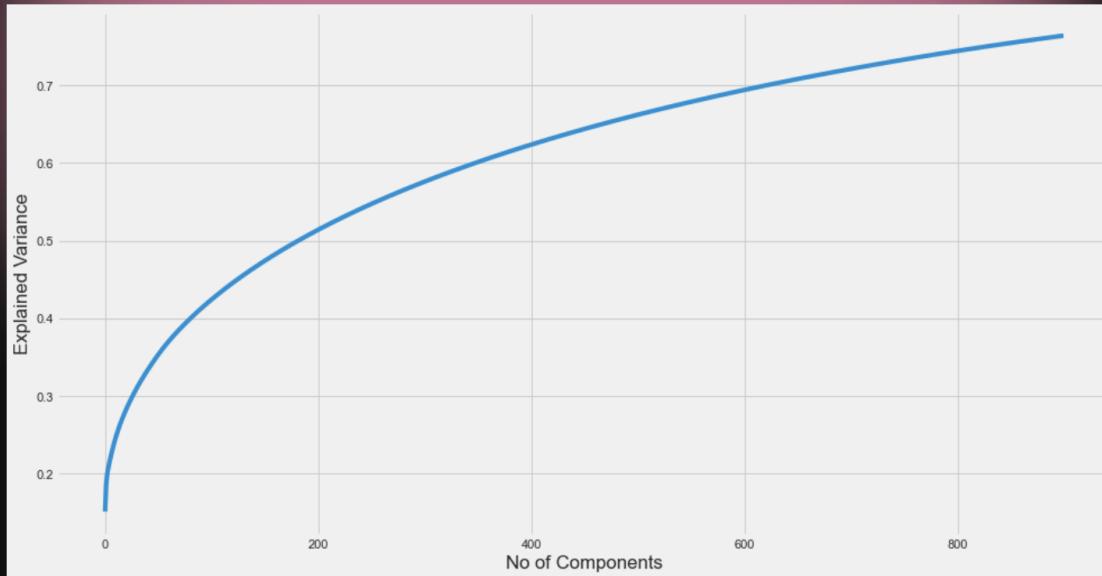
data cleaning

Feature reduction

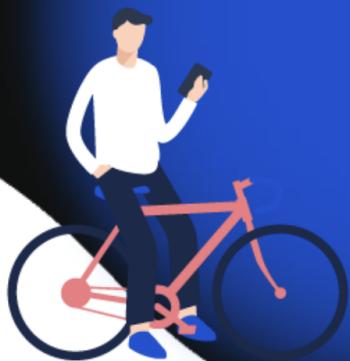
The final dataset

## Feature Reduction

- We applied PCA to reduce the columns to 900.



## Data Preparation

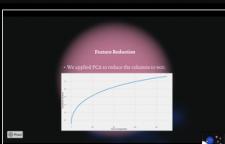


feature  
selection

data cleaning

Feature  
reduction

The final  
dataset



## Final dataset

- 50,000 observation.
- 900 features.



## Data Preparation



feature  
selection

data cleaning

Feature  
reduction

The final  
dataset

Final dataset

\* 10,000 observations

\* 300 features

Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

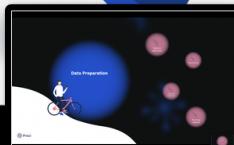
Tools

Modeling

conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl

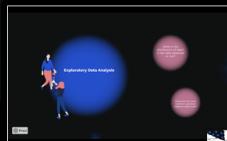
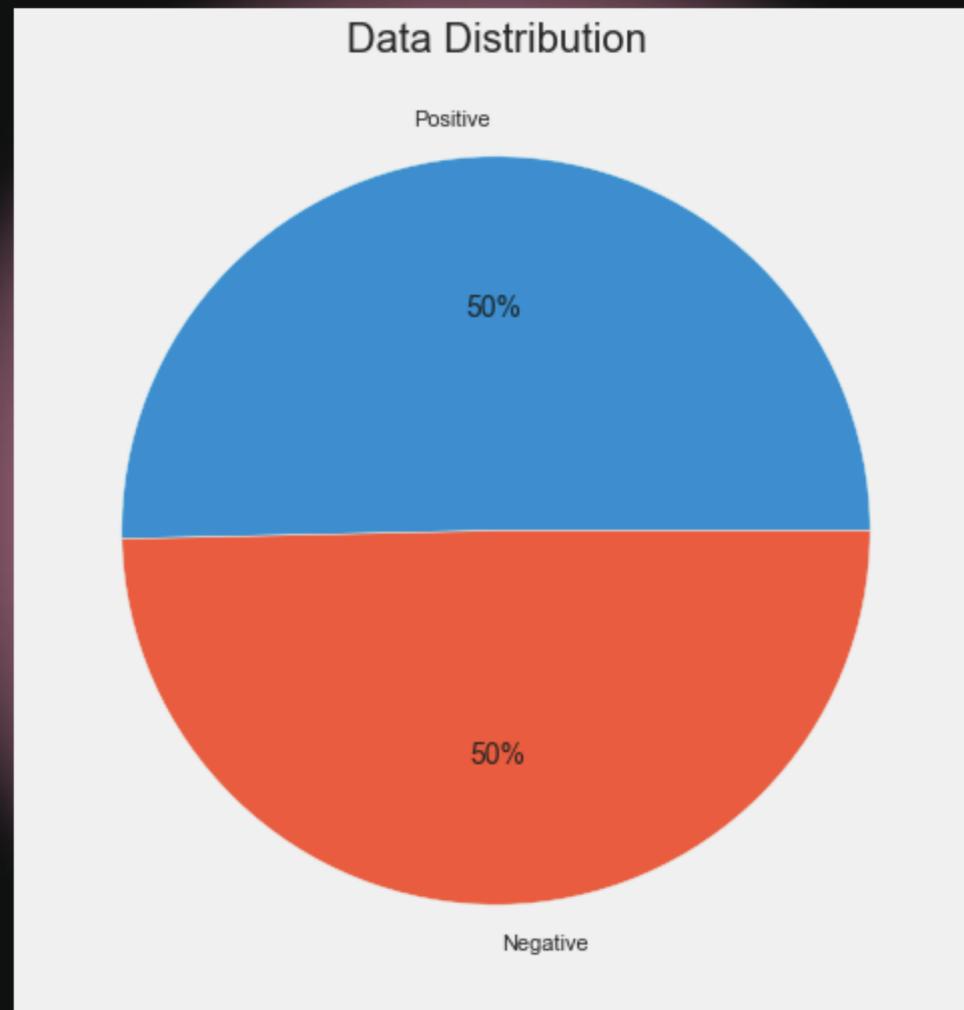




## Exploratory Data Analysis

What is the distribution of data?  
Is the data balanced or not?

What are the most common positive/negative adjectives?

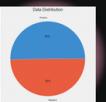


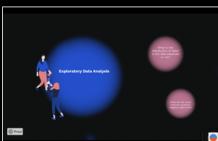
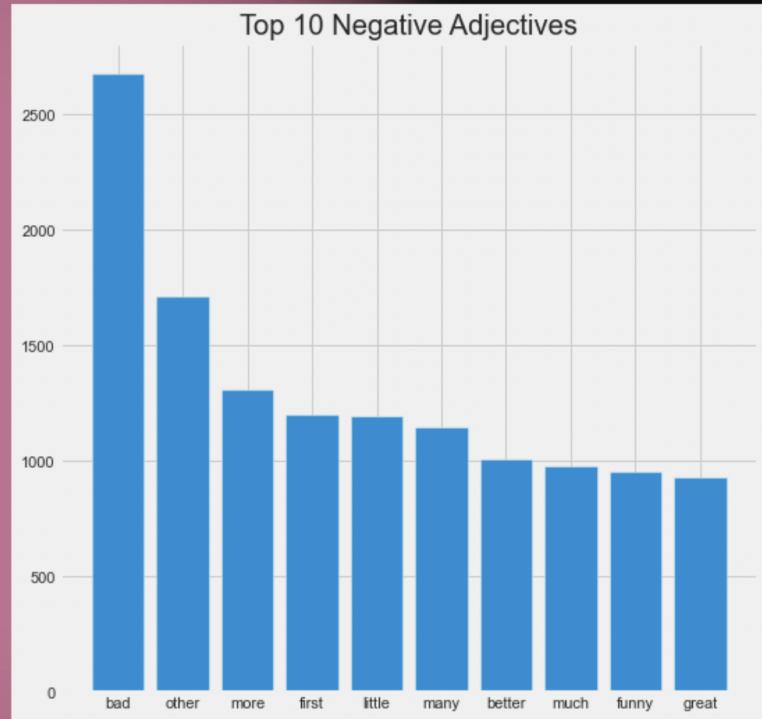
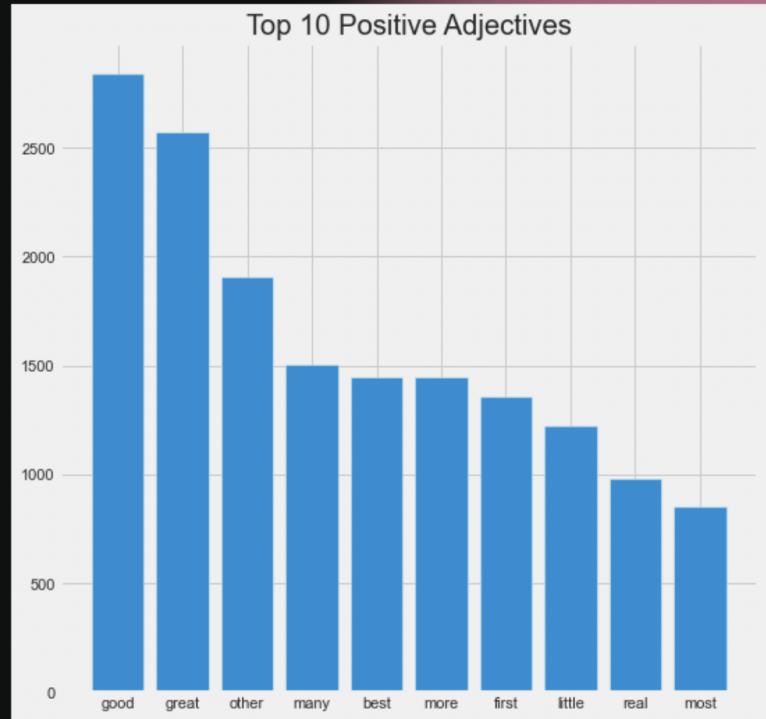


## Exploratory Data Analysis

What is the distribution of data?  
Is the data balanced or not?

What are the most common positive/negative adjectives?



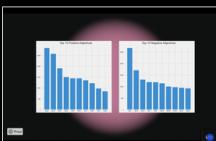




## Exploratory Data Analysis

What is the distribution of data?  
Is the data balanced or not?

What are the most common positive/negative adjectives?



Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

Tools

Modeling

conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl



## Modeling

Models used

Confusion matrix

Result

ROC Curve

- KNN
- Logistic Regression
- Decision Tree
- Random Forest



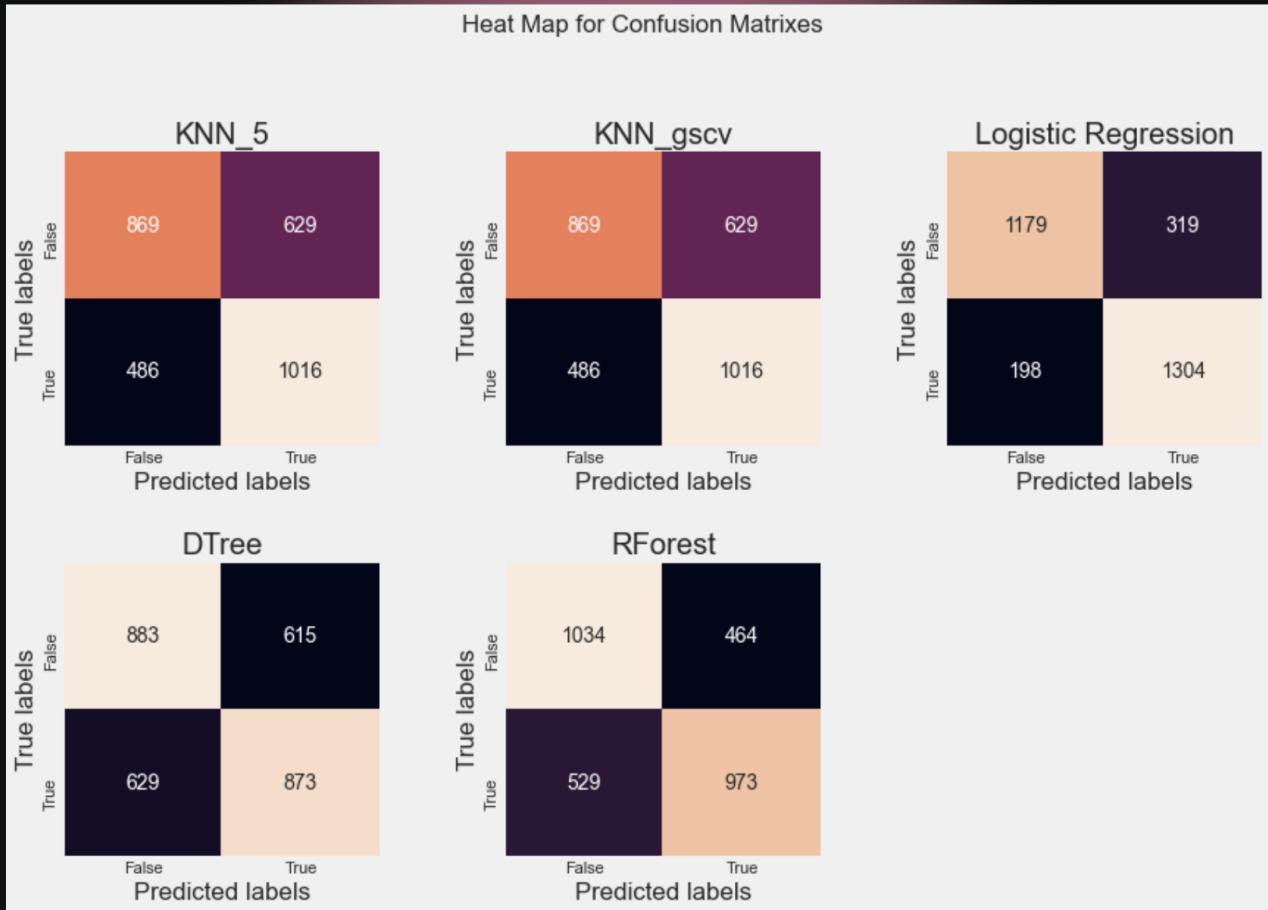
## Modeling

Models used

Confusion matrix

Result

ROC Curve



## Modeling

Models used

Confusion matrix

Result

ROC Curve



### *Models Results*

Model	Accuracy	Precision	Recall	F1
KNN_5	62.83	61.76	67.64	64.57
KNN_gscv	62.83	61.76	67.64	64.57
Logistic Regression	82.77	80.35	86.82	83.46
DTree	58.53	58.67	58.12	58.39
RForest	66.9	67.71	64.78	66.21



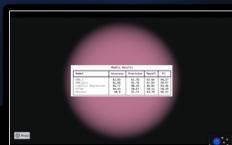
## Modeling

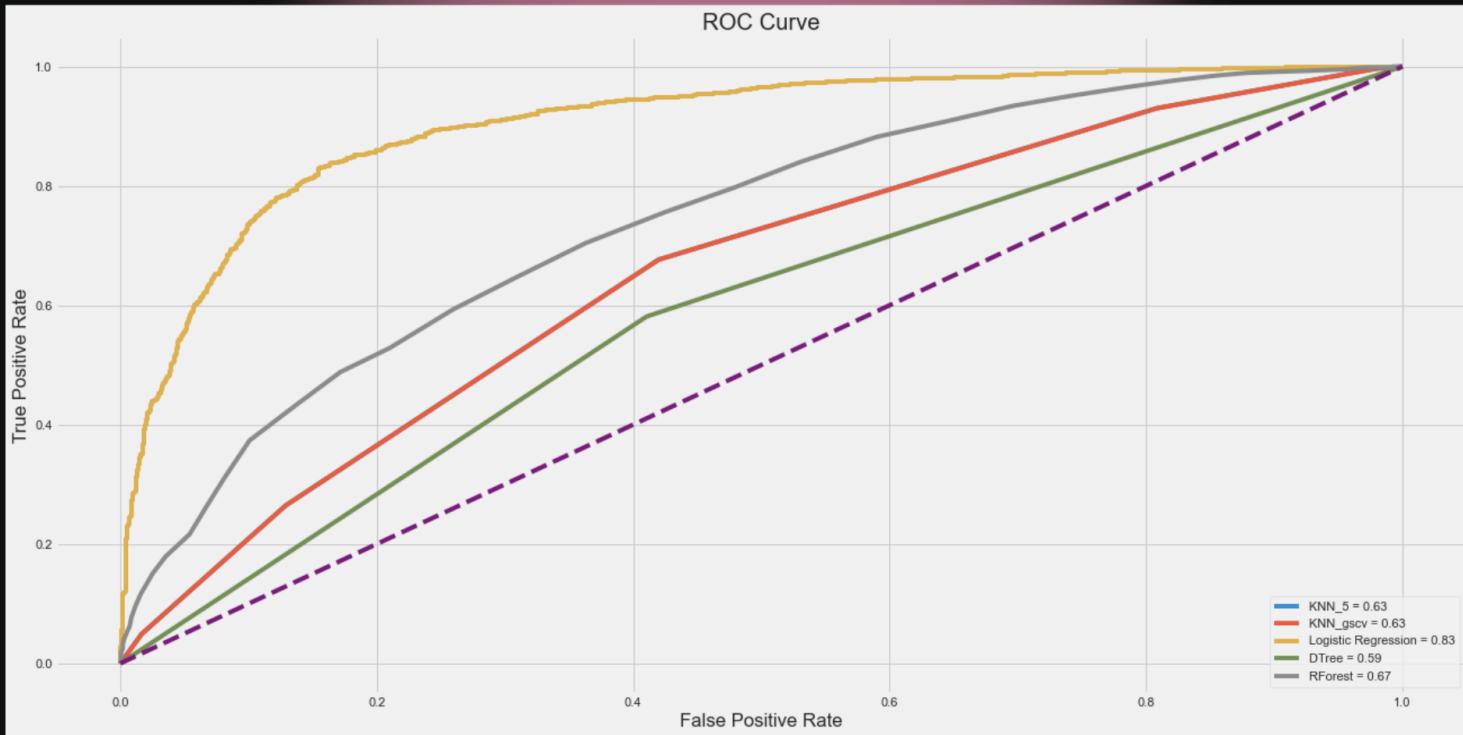
Models used

Confusion matrix

Result

ROC Curve





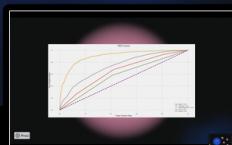
## Modeling

Models used

Confusion matrix

Result

ROC Curve



Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

Tools

Modeling

conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl



## Tools

- Pandas
- Numpy
- Seaborn
- Matplotlib
- Sklearn
- SpaCy



Description  
and  
Objectives

Exploratory  
Data Analysis  
(EDA)

Data  
Preparation

Tools

Modeling

conclusion

# NLP Classification for IMDB Reviews

Faisal Alasgah, Ali Altamimi, Saleh Aljomyl



## Conclusion