

Hate Speech and Offensive Language Recognition

Natural Language Processing

Dimitris Patiniotis Spyropoulos,

July 1, 2022

Table of Contents

- ① Introduction
- ② Methodology
- ③ Results
- ④ Conclusions

Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

Results

Conclusions

Introduction

Hate Speech and Offensive Language Recognition

- Practical solution to filter spam in content sharing apps/websites
- A multi-class classification problem
- Different ways to approach the problem

Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

Results

Conclusions

Methodology

Data Collection and Preparation

Dataset used from Davidson et al. (2017).

[https://www.kaggle.com/datasets/mrmorj/
hate-speech-and-offensive-language-dataset?datasetId=723100&
sortBy=voteCount](https://www.kaggle.com/datasets/mrmorj/hate-speech-and-offensive-language-dataset?datasetId=723100&sortBy=voteCount)

Hate Speech
and Offensive
Language
Recognition

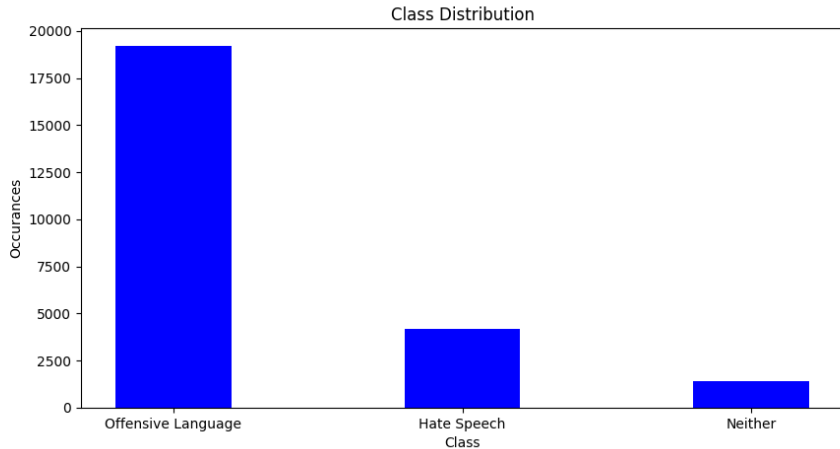
Introduction

Methodology

Results

Conclusions

Class Representation Imbalances



Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

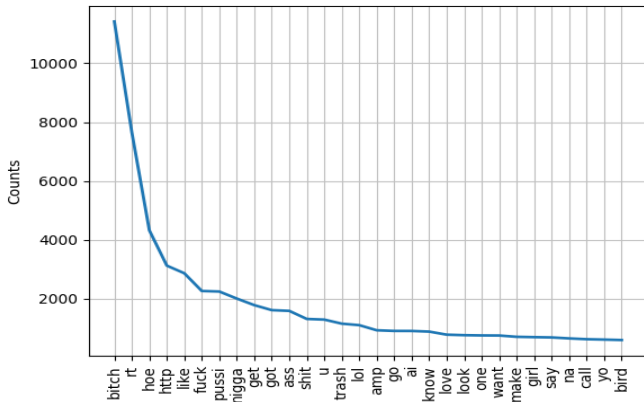
Results

Conclusions

Class Representation Imbalances

- Solved by Random Under-Sampling
- Tested all algorithms to both the original and the balanced dataset

Word Frequency

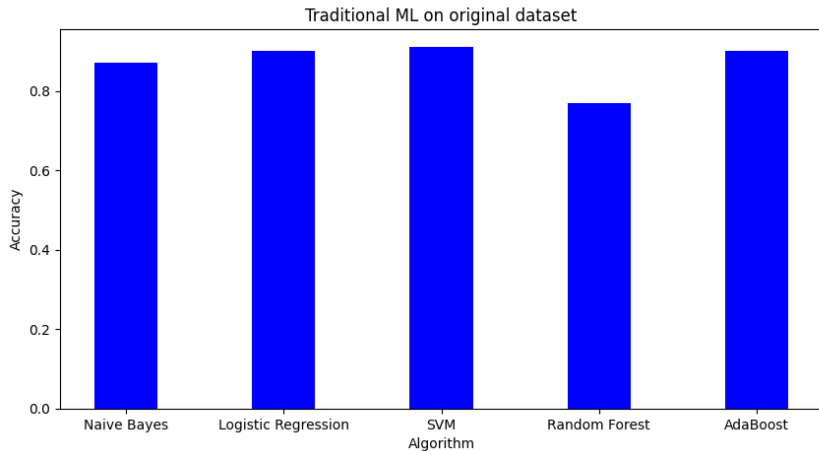


Algorithms Used

- Traditional ML: Naive Bayes, Logistic Regression, SVM, Random Forest, AdaBoost
- LSTMs
- Using Bert

Results

Traditional ML - Original Dataset



Hate Speech
and Offensive
Language
Recognition

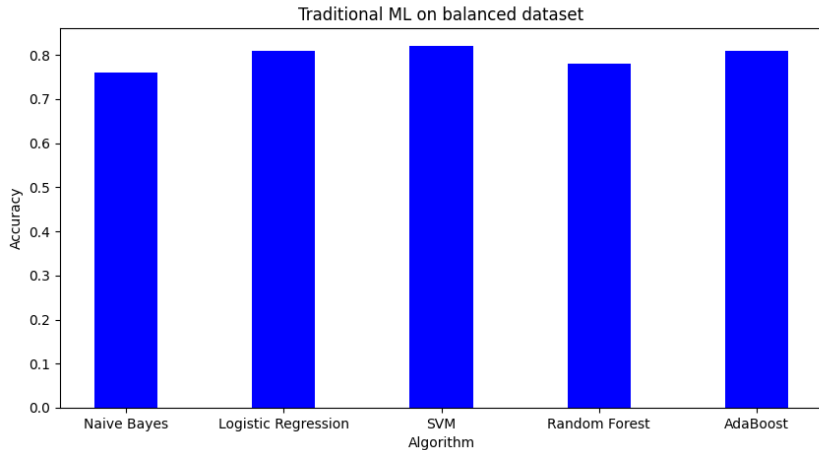
Introduction

Methodology

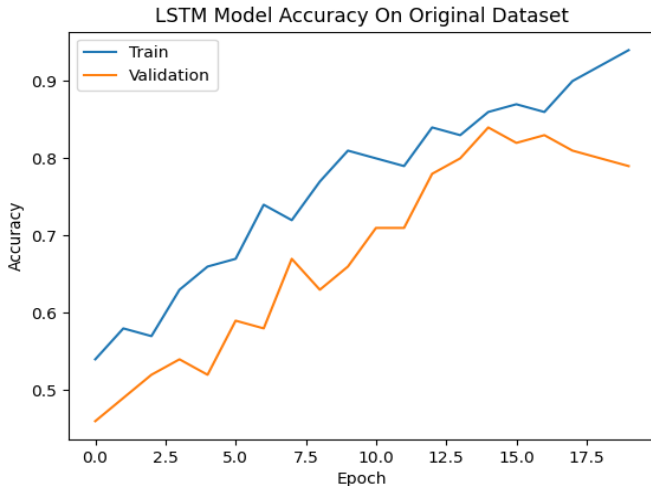
Results

Conclusions

Traditional ML - Balanced Dataset



LSTM - Original Dataset



Hate Speech
and Offensive
Language
Recognition

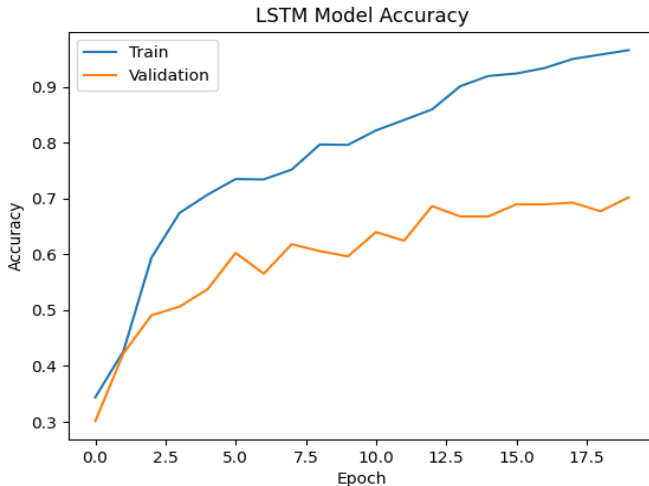
Introduction

Methodology

Results

Conclusions

LSTM - Balanced Dataset



Hate Speech
and Offensive
Language
Recognition

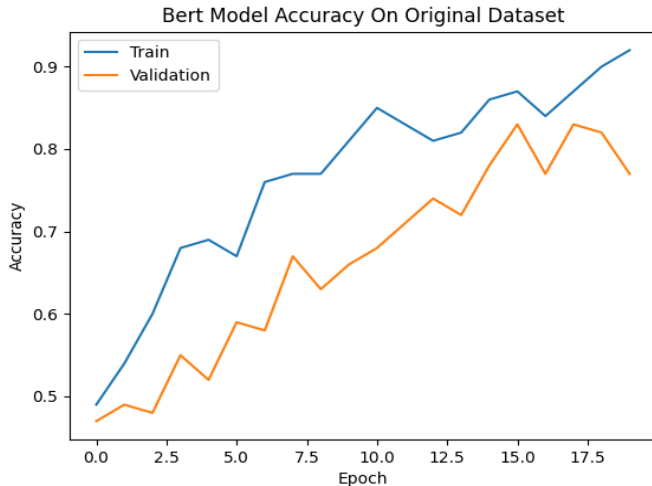
Introduction

Methodology

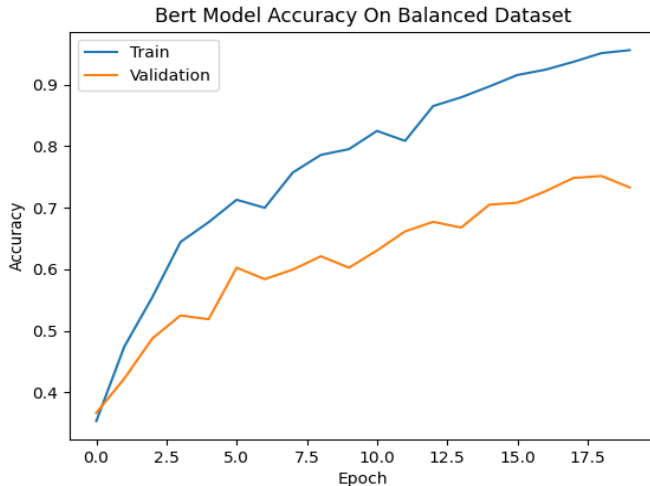
Results

Conclusions

Bert - Original Dataset



Bert - Balanced Dataset



Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

Results

Conclusions

Conclusions

Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

Results

Conclusions

Conclusions

- Larger unbalanced Dataset vs Smaller Balanced Dataset Trade Off
- Traditional ML vs LSTMs
- Introduction of Pretrained Transformer to Improve Performance

Hate Speech
and Offensive
Language
Recognition

Introduction

Methodology

Results

Conclusions