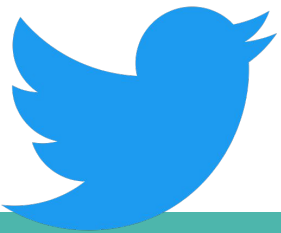

Hate Speech Classifier

— Sirash Phuyal, Sam Macy,
Sushen Kolakaleti —

Introduction

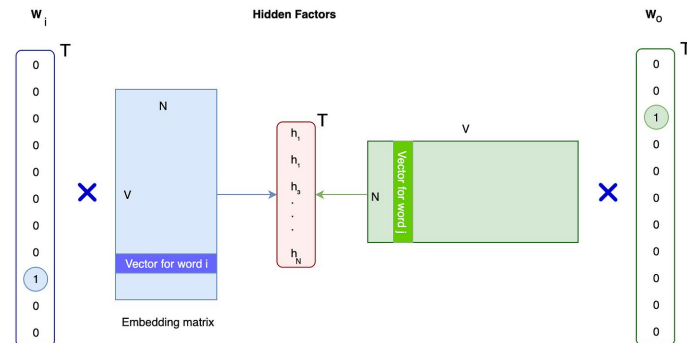
- Hate speech become common social platforms triggers the need for more effective detection mechanisms.
- Project attempts classification of hate speech in tweets, leveraging three distinct embedding approach: BERT-based and GloVe/word2vec
- This threefold embedding exploration aims to offer a comparative analysis on the efficacy of each approach for hate speech detection.



Methodology

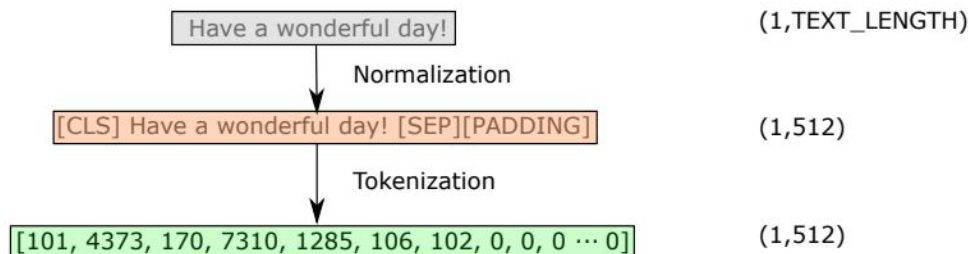
- Dataset processing / cleaning:

- Removing stop-words for GloVe/word2vec
- Majority voting of labels
(Dataset contained multiple labels for many tweets)

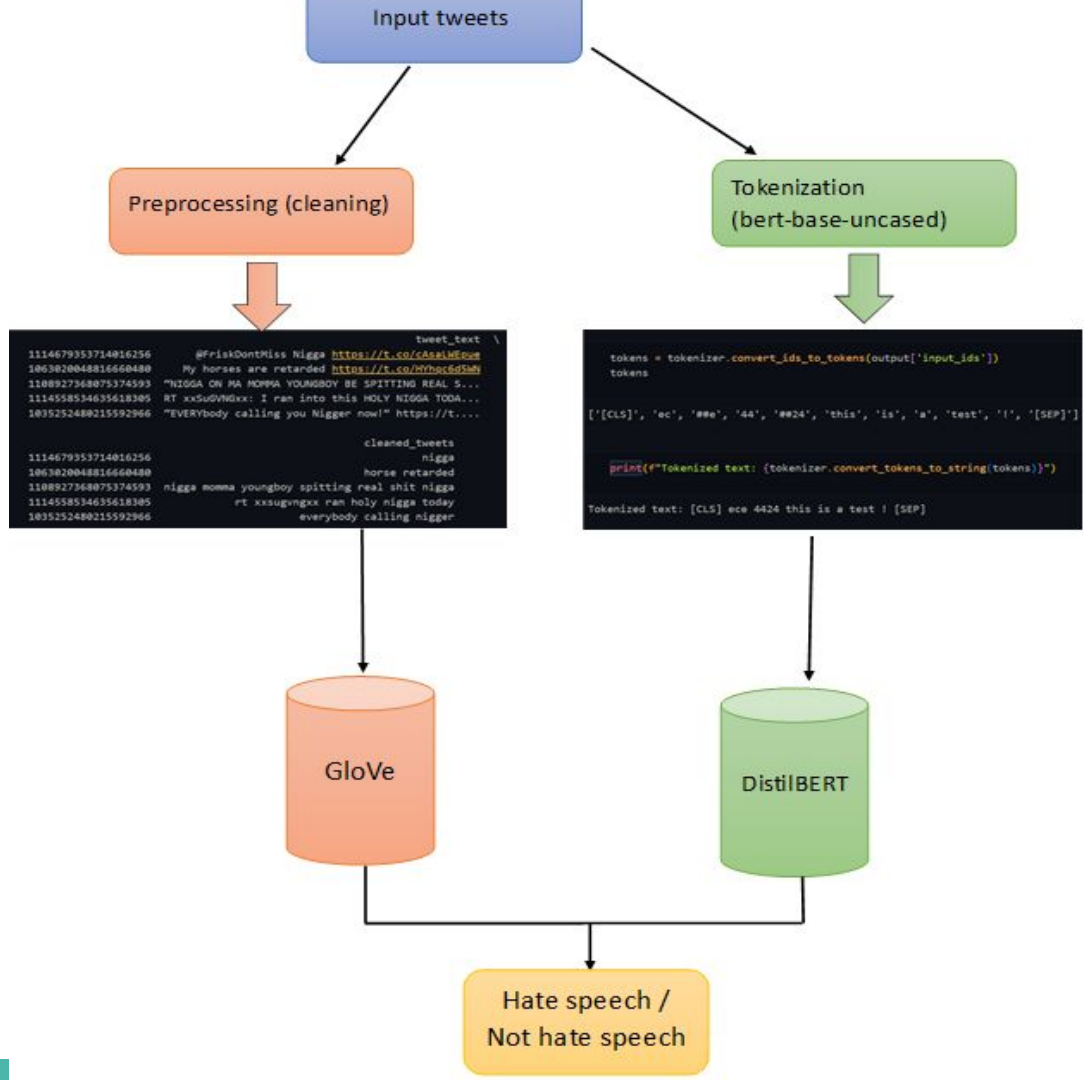


- BERT (Bidirectional Encoders Representation from Transformers)

- Context-rich pretrained model
- Fine-tuned to our dataset, resource-intensive, used condensed version called DistilBERT

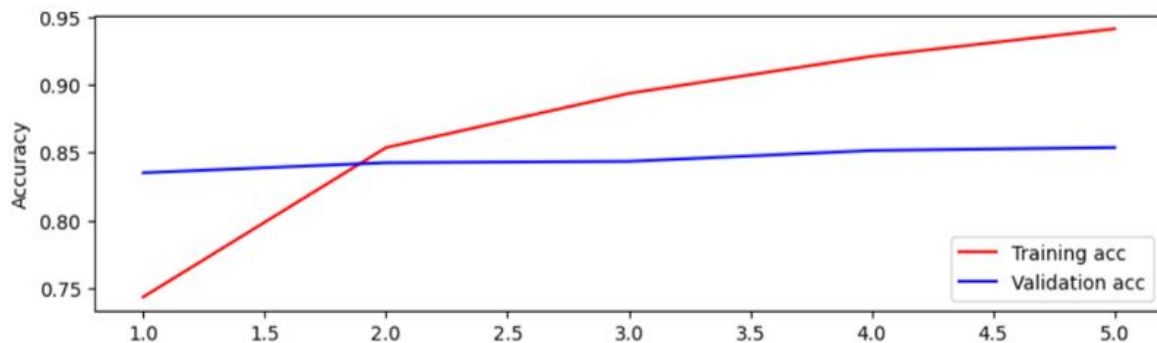


Process



Results & Analysis

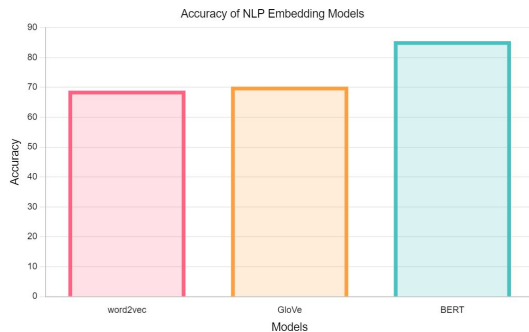
- Word2vec: 73.57% training accuracy, 68.92% testing accuracy
- GloVe: 69.71% training accuracy, 70.24% testing accuracy
- BERT: 45.34% loss, 85.48% testing accuracy



Training vs validation accuracy: Validation consistent around 85%

Conclusion

- Performance is highly dependent on the quality of the dataset.
- BERT performance and accuracy better than GloVe/word2vec
- BERT initially applied stop-word removal and lemmatization before tokenization, it is context-aware
- Potentially even better performance with a larger scale Deep NN model



Any Questions

