

Hate detection by NLP and deep learning methods

...

uc3m

Universidad
Carlos III
de Madrid

Presented By: Rami Abulfadl
Advisor: Isabel Segura-Bedmar
Masters in Big Data | September 2020

Content

- Definition and Motivation
- Objectives
- Approach: dataset and methods
- Evaluation
- Conclusions and Future work

What is hate?

“any kind of communication that attacks or uses discriminatory language to a person or a group on the basis of who they are, in other words, based on their religion, ethnicity, nationality, race, colour, descent, gender or other identity factor.”



UNITED NATIONS STRATEGY AND PLAN OF ACTION ON HATE SPEECH

why is it important to detect ?

**Facebook engineers quit in protest, accusing the company
of 'profiting off hate'**



“When the looting starts, the shooting starts” Donald J. Trump

why is it important to detect ?

- Remove abusive content for a good cause.
- Business reputation



Donald J. Trump  @realDonaldTrump · 2h

This Tweet violated the Twitter Rules about glorifying violence. However, Twitter has determined that it may be in the public's interest for the Tweet to remain accessible. [Learn more](#)

View

A screenshot of a Twitter post from Donald J. Trump (@realDonaldTrump). The post was made 2 hours ago and contains a warning message. The message states: "This Tweet violated the Twitter Rules about glorifying violence. However, Twitter has determined that it may be in the public's interest for the Tweet to remain accessible." It includes a link to "Learn more". The "View" button is located at the bottom right of the message box. The background of the slide is dark blue.

Objectives

Acquire

the basic knowledge
in the field of NLP
and deep learning

Analyse

the progression of
methods and
techniques till
reaching the state of
art performance in
hate detection.

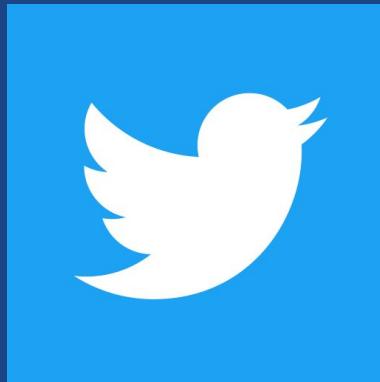
Architect

two most advanced
deep learning models,
BERT and Bi-LSTM,
for hate detection
using Twitter dataset.

Assess

The performance of
models through error
analysis.

Approach: Twitter dataset



32k tweets

Train: 80%

Validation: 10%

Test: 10%

Neutral

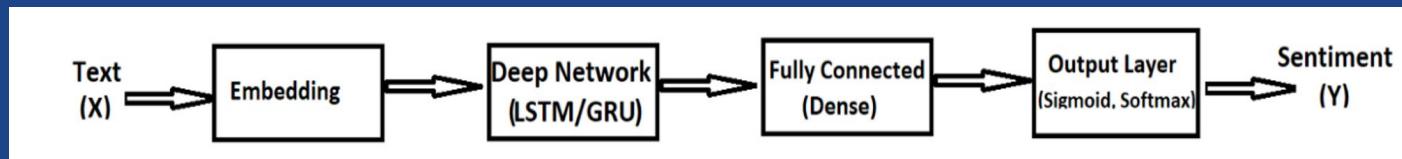
29,720

Hate

2,242

First Approach

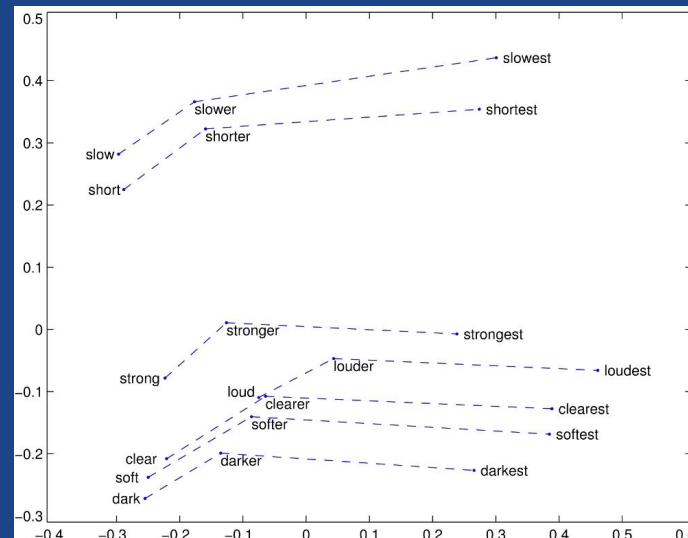
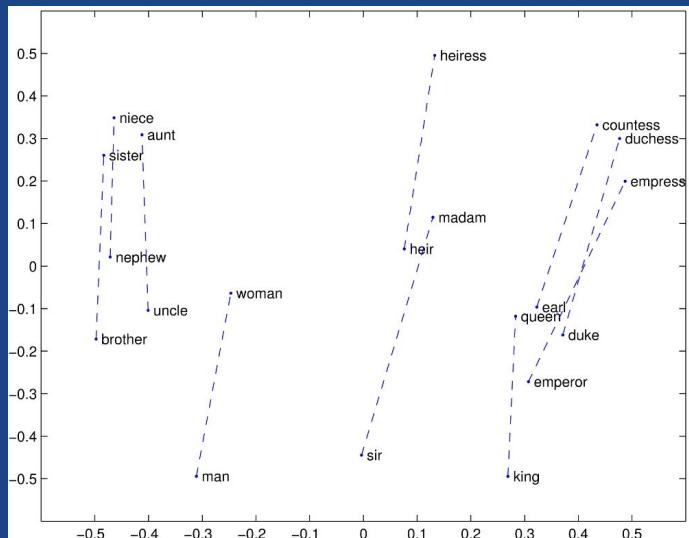
Bi-LSTM



Bi-LSTM

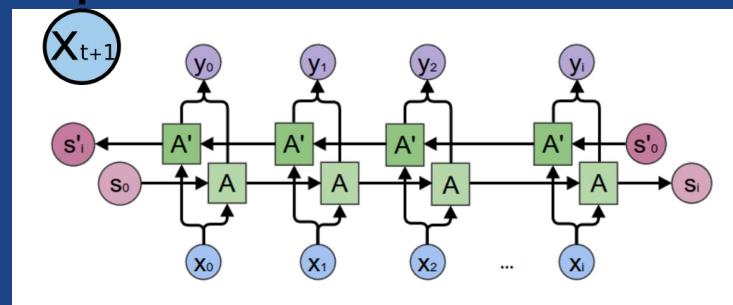
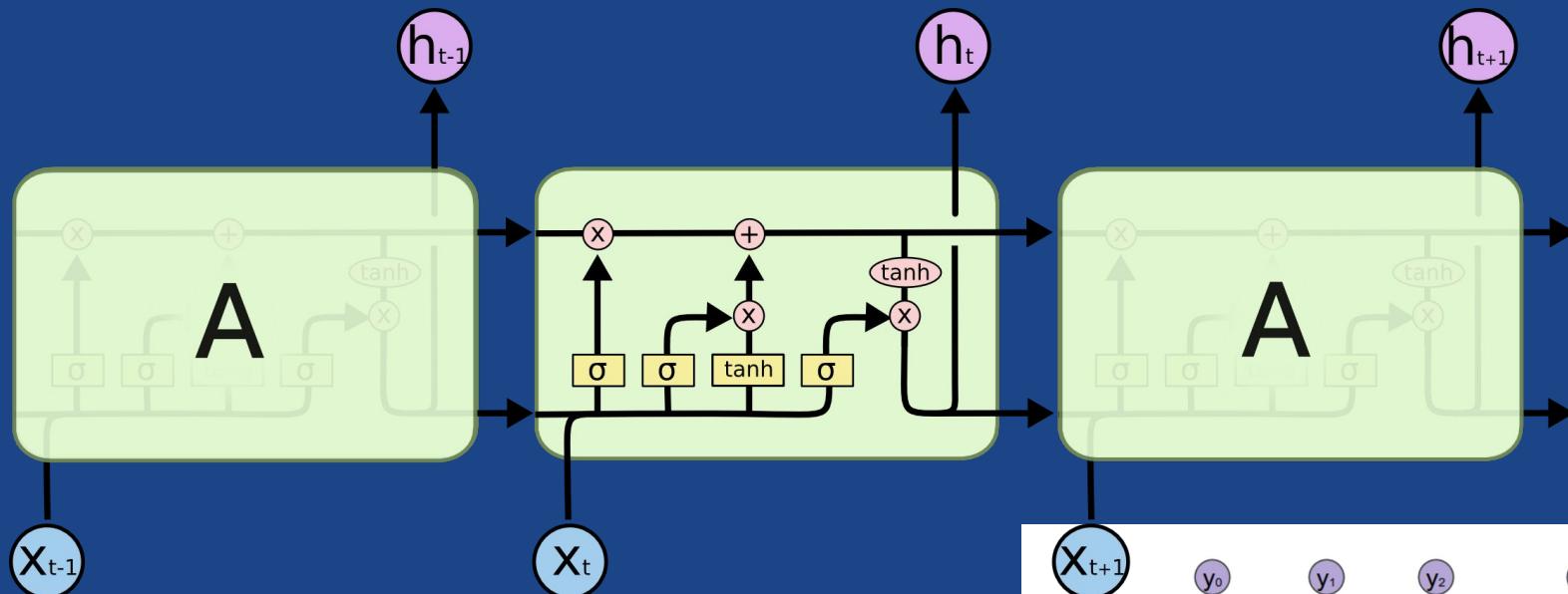
Word Embedding by GloVe for vector space representations

glove.6B.100d 6 B tokens of 400 k vocab forming a 100 dimension form Wikipedia 2014 and Gigaword 5 corpora



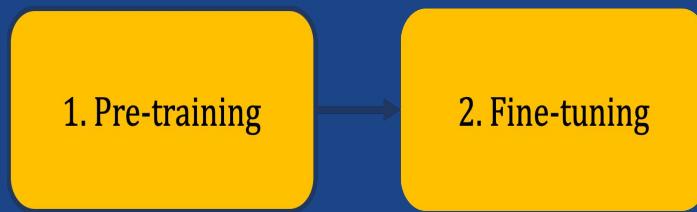
Pennington, J., Socher, R., & Manning, C. D. (2014, October). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)* (pp. 1532-1543).

Bi-LSTM

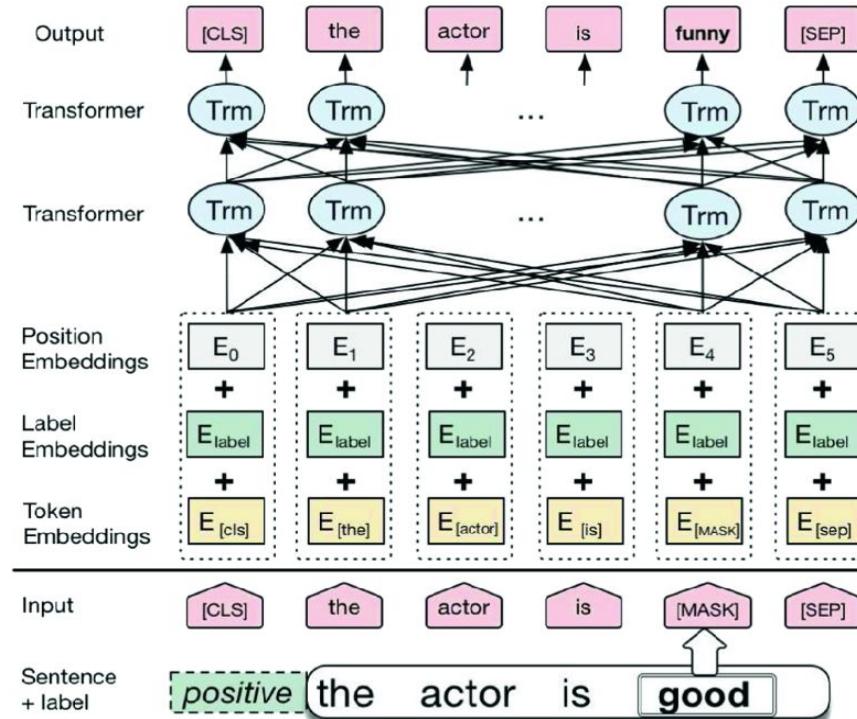


Second Approach

BERT

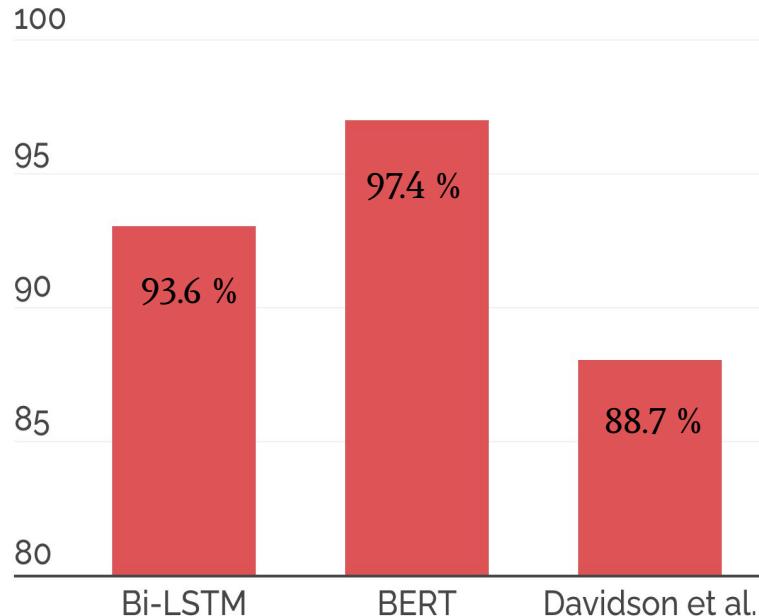


BERT

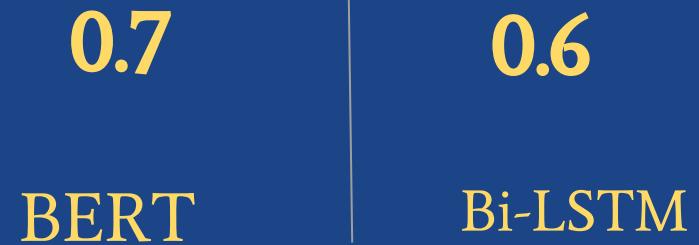


Evaluation

Accuracy %

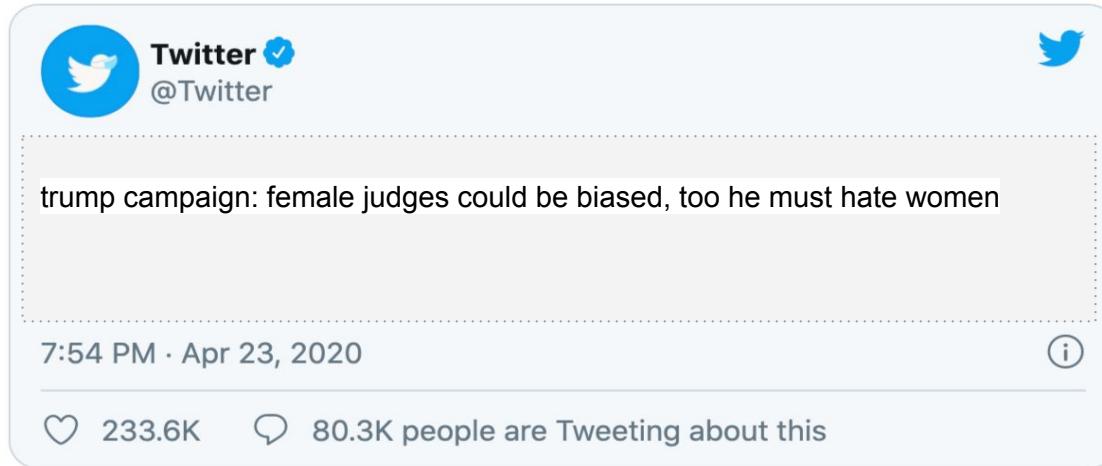


F1-Score



Error Analysis

False positive by LSTM, but True negative by BERT



Error Analysis

False negatives by BERT but True positive by LSTM

Twitter

@Twitter

al capone was jailed for tax evasion, not murders he committed. what is
@user's achilles heel? #unpresidented #nohate

7:54 PM · Apr 23, 2020

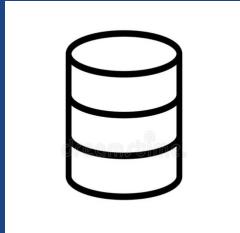
233.6K 80.3K people are Tweeting about this

(i)

Conclusions

- BERT achieve the best performance in classification.
- Both Models get higher accuracy than original davidson et al. results.
- Both Models were able to handle unbalanced classes better.

Future Work



Big dataset for
training.



Develop better
word Embedding
for slang words
and
abbreviations



Cross lingual unified
model for non
english languages

Code repository on GitHub

<https://github.com/ramiabulfadl/Hate-detection>

Thank
you!