

WHO • SAID • THAT • ?

Comparing performance of TF-IDF and fastText to identify authorship of short sentences

PROBLEM

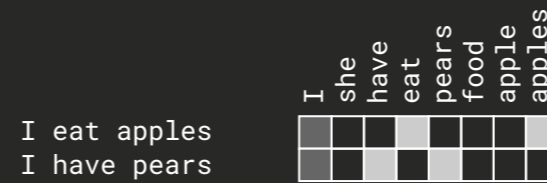
Identifying authorship of text can be done by matching features extracted from text.

Currently mostly performed on large text documents, not on short everyday sentences, which could benefit personal assistants or chatbots.

Compare performance of two popular extraction techniques **TF-IDF** and **fastText** by answering the questions:

- Which mistakes does the model make?
- Does performance change as sentence length increases?

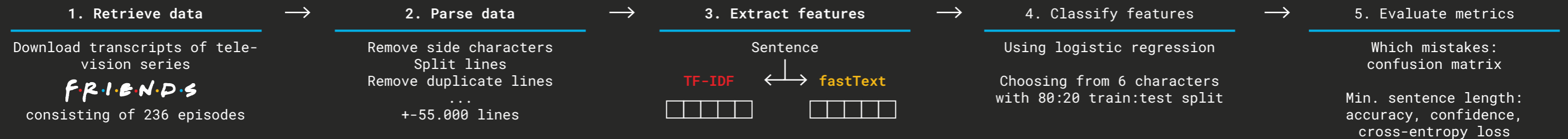
TF-IDF measures importance of words by counting occurrences in document



fastText represents words with similar meanings and context the same way



METHOD



RESULTS

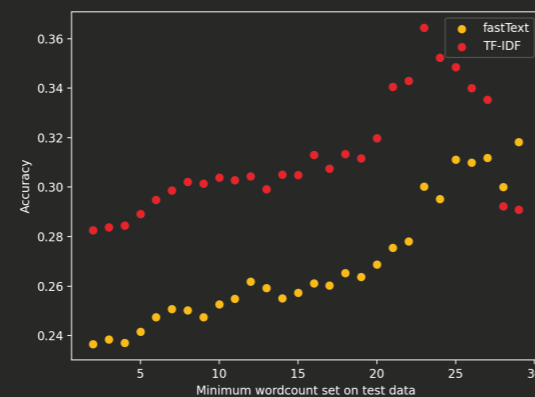
Confusion matrix

Amount of correct predictions

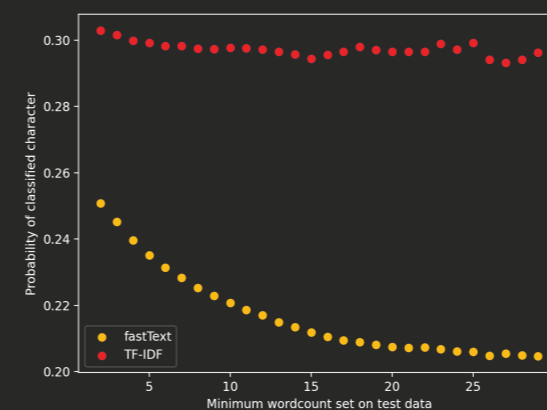


fastText	332	294	569	673	331	515
TF-IDF	457	408	656	636	433	654

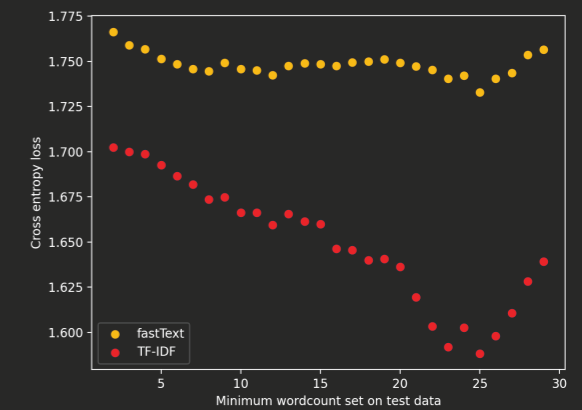
Accuracy



Confidence



Cross-entropy loss



CONCLUSION

TF-IDF outperforms **fastText** in every measurement, but its performance is only slightly better than randomly guessing the original character, reaching an accuracy of 28 percent when making a distinction between 6 characters.

Accuracy increases linearly at the same rate for both techniques test set's sentence length increases.

TFIDF's confidence remains constant as this limit is set on either the test or training data, whereas fastText's confidence decreases and increases, respectively.

Cross-entropy loss, however, remains constant for fastText and decreases for TF-IDF as the minimum word count set on the test data increases.



Thomas van Tussenbroek
CSE3000, 25/06/2020

David Tax; Marco Loog; Tom Viering; Arman Naseri Jahfari; Stavros Makrodimitis;