

Plaint-or-Tweet

Tommaso Battistini, Edoardo De Matteis, Leonardo Emili,
Mirko Giacchini, Alessio Luciani

December 24, 2020

Introduction

Our dataset is [3]

Model

0.1 Classic Naive Bayes

0.2 Embeddings

Preprocessing

Tests and metrics

...

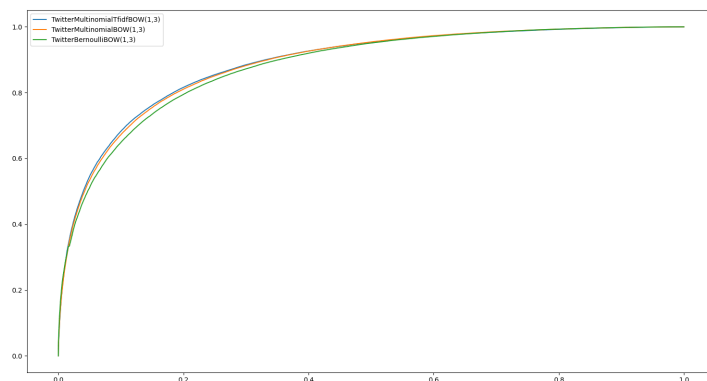


Figure 1: ROC curve for the Naïve Bayes approach with tf-idf.

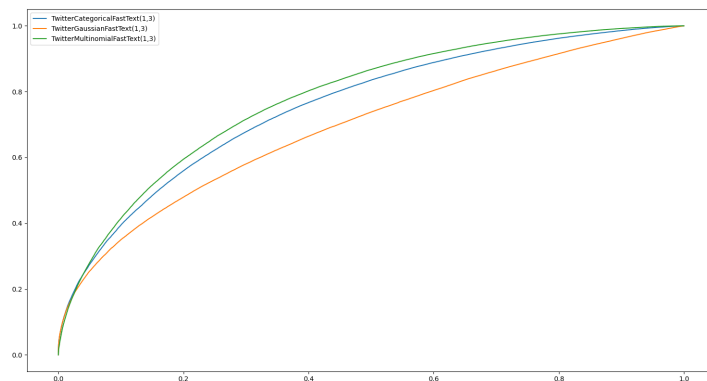


Figure 2: ROC curve for the FastText embedding.

Kaggle notebooks

Further tests

We trained our model with a dataset and tested it against a different one, to this aim were used an IMDb dataset of cinematographic reviews [2] and a Reddit

one about various NFL games [1], results for the Bag-Of-Words Naïve Bayes (1,3)-gram model can be seen in figures 3 and 4.

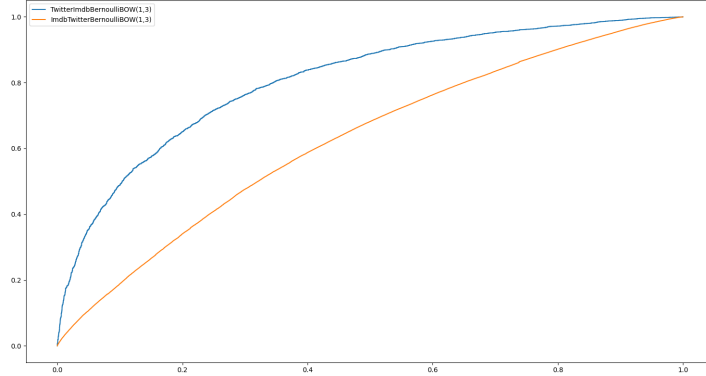


Figure 3: Comparison between Imdb-Twitter and Twitter-Imdb.

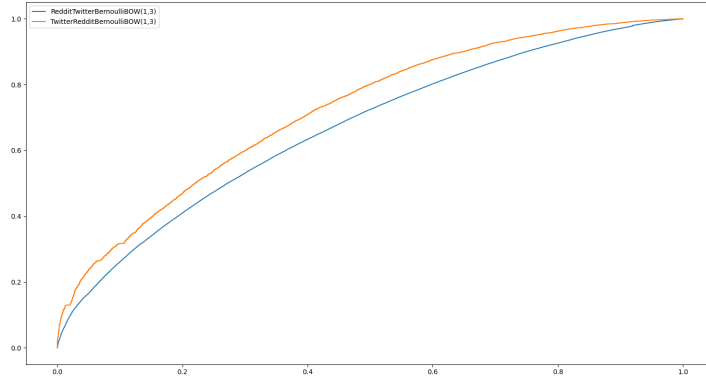


Figure 4: Comparison between Reddit-Twitter and Twitter-Reddit.

This is something usually not done since different datasets will have different distributions hence we expect bad results, strangely enough training our model with the Twitter dataset and testing it with the IMDB one gives rather good metrics that can be seen in table 1, is exposed only the comparison between Twitter and IMDB since results are more interesting. This could happen since Twitter’s dataset has a very large number of examples, other datasets other than being smaller are also very specific: the Reddit one is only about a specific

set of american football games and the IMDb one is about cinematographic reviews, reasonably in this case lexicon is also more specific. Sentiment140 on the other hand covers a larger class of human language and is less prone to be biased.

Table 1: Comparing metrics for some train test combinations.

Train-Test	Accuracy	F1	AUROC
Twitter-IMDb	0.732	0.723	0.806
IMDb-IMDb	0.891	0.887	0.963
IMDb-Twitter	0.550	0.330	0.625
Twitter-Twitter	0.797	0.800	0.880

References

- [1] Caio Brighenti. Nfl draft reddit comments. URL: <https://www.kaggle.com/caiobrightenti/nfl-draft-reddit-comments?select=picks.csv>.
- [2] IMDb. Imdb datasets. URL: <https://www.imdb.com/interfaces/>.
- [3] Marios Mikaelides Kazanova. Sentiment140, 2017. URL: <https://www.kaggle.com/kazanov/sentiment140>.