

Unsupervised Embeddings with Graph Auto-Encoders for Multi-domain and Multilingual Hate Speech Detection

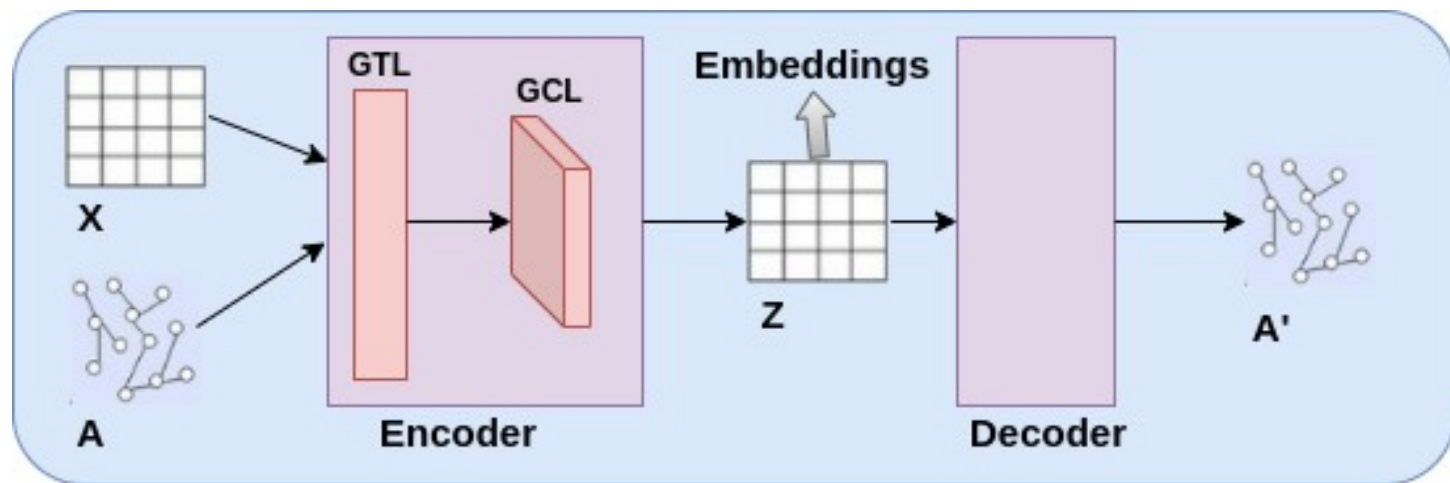
Gretel Liz De la Peña Sarracén^{1,2}
Paolo Rosso¹

¹Universitat Politècnica de València

²Symanto Research

LREC, June 2022

Graph Auto-Encoders



Our model - GAE

XHate-999 dataset

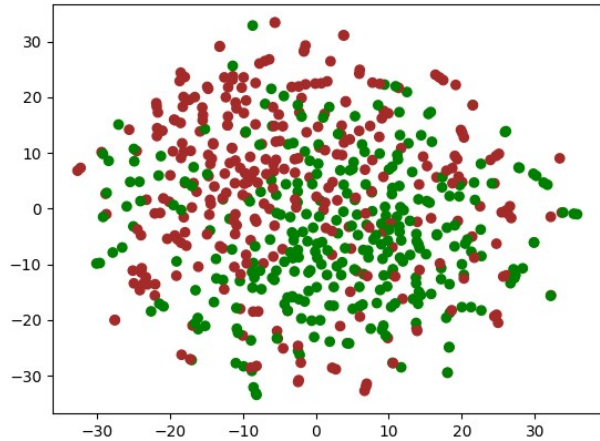
6 languages:

- English (EN)
- German (DE)
- Russian (RU)
- Turkish (TR)
- Croatian (HR)
- Albanian (SQ)

For each language there are 3 domains:

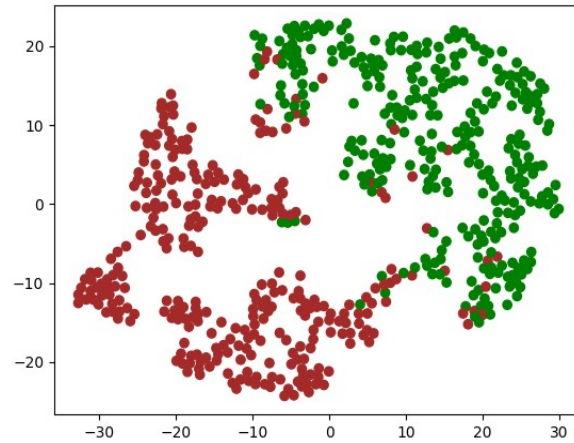
- Fox News (Gao) with 99 samples
- Twitter/Facebook (Trac) with 300 samples
- Wikipedia (Wul) with 600 samples

Analysis of Latent Representation - Embeddings

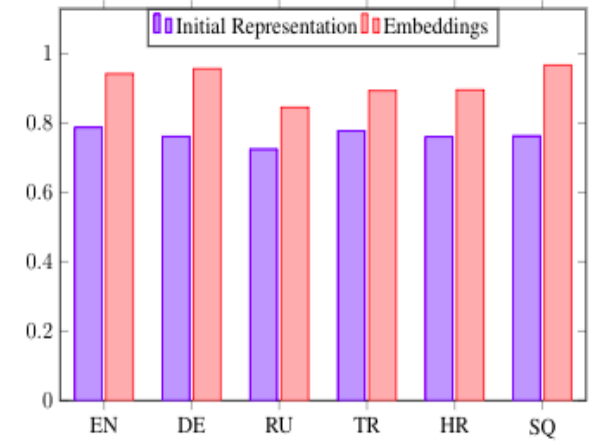


Initial representation

English - Wul

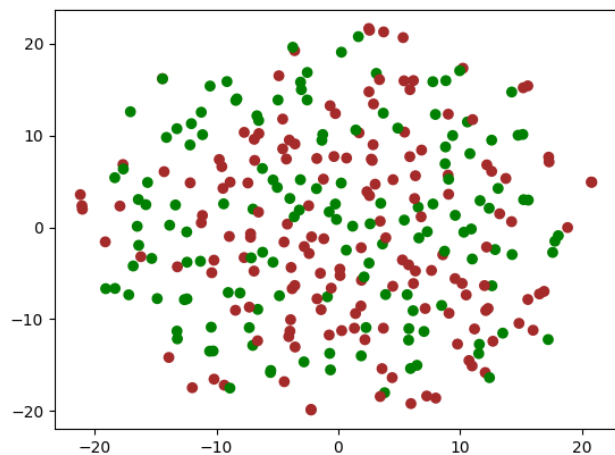


Embeddings



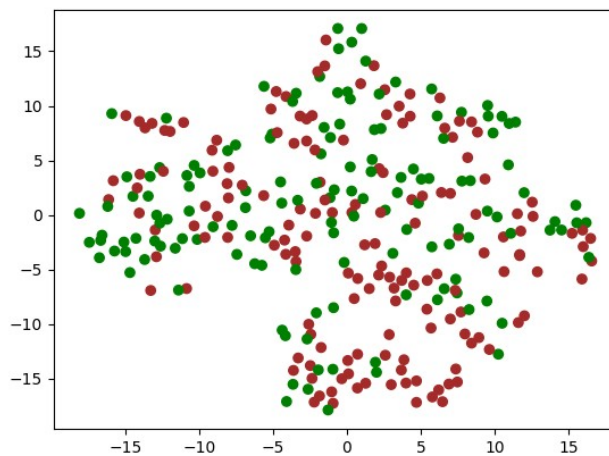
F1 in Wul.
Hate Speech Detection

Analysis of Latent Representation - Embeddings

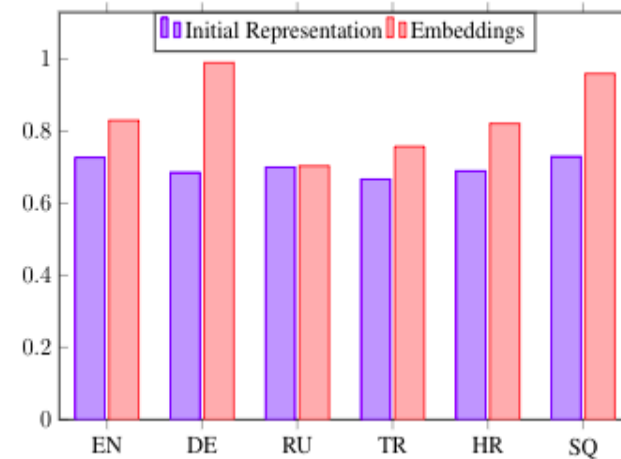


Initial representation

Russian - Trac



Embeddings



F1 in Trac.
Hate Speech Detection

Evaluation for Hate Speech Detection

	GAE		
	Gao	Trac	Wul
EN	0.9714 _{.018}	0.8301 _{.025}	0.9424 _{.007}
DE	0.9565 _{.022}	0.9900 _{.017}	0.9569 _{.011}
RU	0.6667 _{.014}	0.7238 _{.010}	0.8453 _{.014}
TR	0.7826 _{.008}	0.7567 _{.030}	0.8936 _{.021}
HR	0.8799 _{.030}	0.8214 _{.023}	0.8958 _{.017}
SQ	0.9565 _{.011}	0.9600 _{.011}	0.9673 _{.026}

	XLM-R		
	Gao	Trac	Wul
EN	0.5909 _{.041}	0.7245 _{.022}	0.8538 _{.034}
DE	0.6857 _{.031}	0.7272 _{.014}	0.8635 _{.049}
RU	0.5754 _{.009}	0.7070 _{.005}	0.8384 _{.041}
TR	0.5171 _{.011}	0.7371 _{.017}	0.8199 _{.036}
HR	0.5050 _{.047}	0.6377 _{.041}	0.8384 _{.047}
SQ	0.5642 _{.041}	0.7148 _{.016}	0.8231 _{.041}

F1 in Multi-domain Hate Speech Detection

	All		
	GAE	XLM-R	mBERT
EN	0.8333 _{.010}	0.5642 _{.047}	0.5111 _{.053}
DE	0.9565 _{.014}	0.4545 _{.038}	0.4850 _{.045}
RU	0.8333 _{.001}	0.4923 _{.061}	0.4527 _{.031}
TR	0.8799 _{.004}	0.6192 _{.043}	0.3864 _{.057}
HR	0.8461 _{.012}	0.6459 _{.039}	0.4545 _{.044}
SQ	0.9565 _{.002}	0.4978 _{.021}	0.4457 _{.046}

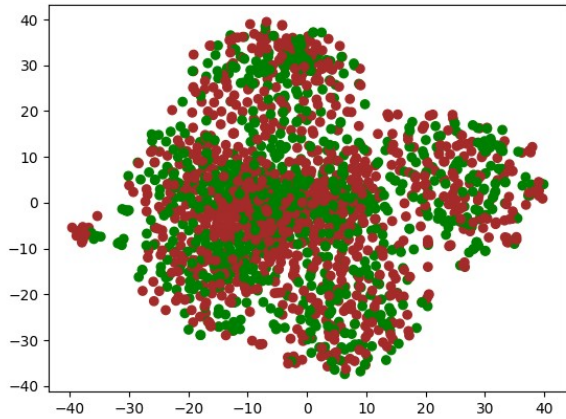
F1 in Multi-domain Hate Speech Detection.
(All domains to train)

Multilingual Evaluation

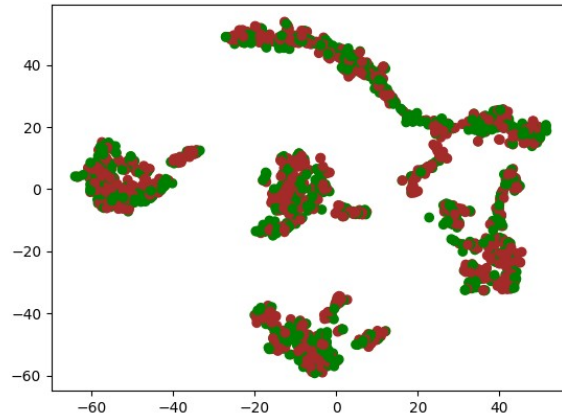
Domain	Gao	Trac	Wul
GAE	0.3972 _{.090}	0.6858 _{.062}	0.6255 _{.058}
GAE-USE	0.9308 _{.011}	0.9598 _{.005}	0.9491 _{.021}
mBERT	0.7047 _{.054}	0.7952 _{.080}	0.8939 _{.072}
XLM-R	0.7349 _{.015}	0.8585 _{.012}	0.9303 _{.008}

F1 in Multilingual Hate Speech Detection

Multilingual Evaluation

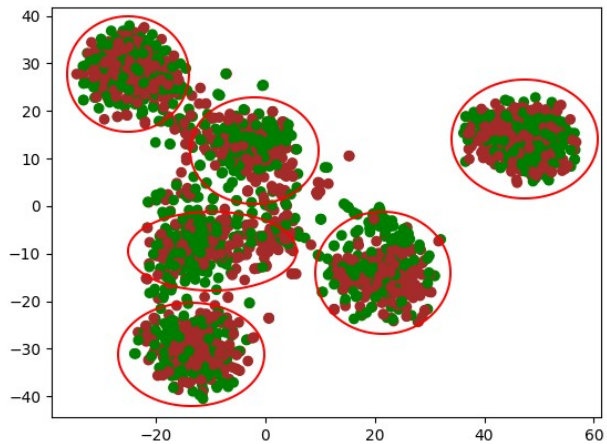


Initial representation
with TFIDF

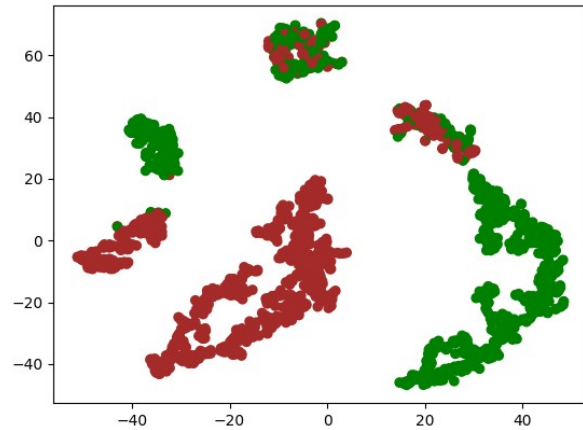


Embeddings

Multilingual Evaluation



Initial representation with
Universal Sentences Encoder



Embeddings

Conclusion

- We generated promising embeddings to discriminate between hateful and non-hateful texts, by using the graph auto-encoder model.
- The model, based on the generated embeddings, achieved good results compared to mBERT and XLM-R.

Unsupervised Embeddings with Graph Auto-Encoders for Multi-domain and Multilingual Hate Speech Detection

Gretel Liz De la Peña Sarracén^{1,2}
Paolo Rosso¹

¹Universitat Politècnica de València

²Symanto Research

LREC, June 2022