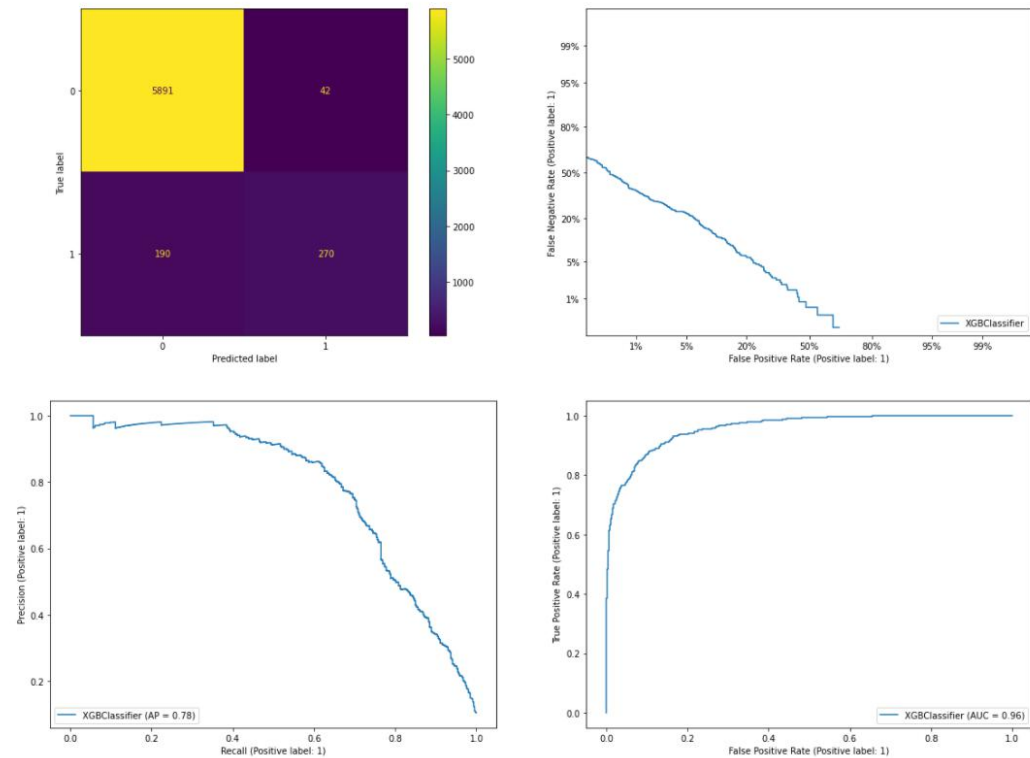# Twitter sentiment analysis

SOUNAK MANDAL

This project attempts at classifying tweets into racist and non-racist categories. The model figures out the connotations behind the word and how likely it is to be linked to racism. The words are encoded by different word embedding techniques. The statistical or frequency based embedding techniques include count vectorization and TF-IDF. The neural word embedding is using word2vec model and is simply averaged to obtain encoding for each tweet.

The models used were standard logistic regression, support vector machines, random forest and xgboost. With all models neural embedding worked best while statistical embedding performed poorly. Finally the xgboost model was fine tuned.
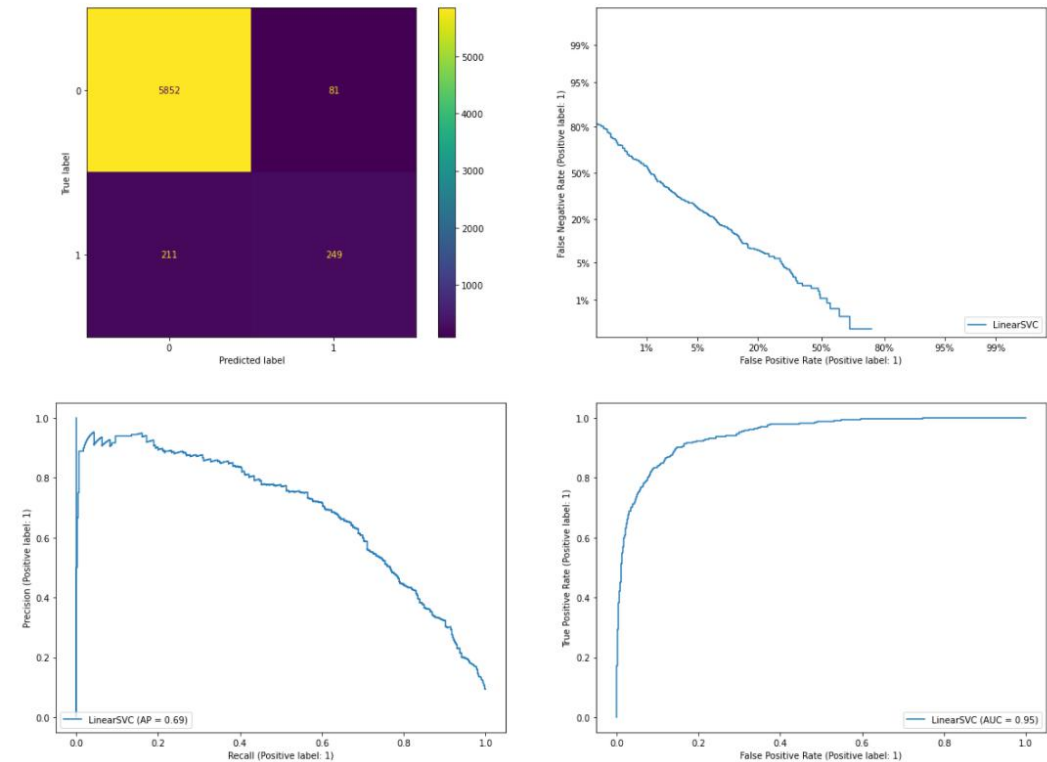
GitHub link : ML-Projects/twitter sentiment analysis at main · SounakMandal/ML-Projects (github.com)

Some of the plots available in the notebook. For each embedding, model combination confusion matrix, det curve, precision recall curve and roc curve were plotted.

# Thank You