

# A Tweet Text Binary Artificial Neural Network Classifier

Theodore Nikolettopoulos<sup>1</sup>, Claudia Wolff<sup>2</sup>

<sup>1</sup>Unaffiliated, Athens, Greece

<sup>2</sup>Kiel University, Kiel, Germany

theo\_nikolettopoulos@yahoo.co.uk, wolff@geographie.uni-kiel.de

## ABSTRACT

We present an Artificial Neural Network (ANN) text classifier to deal with the task of automatically detecting a tweet as being flood-related or not. The framework for classifying flood-related tweets consists of three basic ANN models. Each model is a different ANN type and the final output is determined by a majority rule on the individual model outputs. The overall F1 score on the test set was 0.5405, significantly lower than on the training/validation set, suggesting that we overfitted the training set.

## 1 INTRODUCTION

This research was conducted as part of the ‘Flood-Related Multimedia Task’ challenge provided by the Multimedia Evaluation Benchmark (MediaEval) 2020 [1]. The goal of the task is to automatically identify and classify tweets which are relevant to flooding in Northeastern Italy. For this binary classification problem, we used different types of ANNs to automatically classify the tweet’s text [2]. As different types of ANNs might capture different characteristics of the ANN input, we chose to implement three different types and determine the final decision by using a majority rule on the individual ANN outputs.

## 2 APPROACH

### 2.1 Text Vectorization

To convert the tweet’s text to a numeric format as required by the ANNs input layers we make use of word embeddings [3]. Word embeddings are a way to map words onto low dimensional (compared to other text numerical representation formats) vectors with the important property that words with similar meaning are mapped to vectors which are close to each other (in e.g. Euclidean distance) in the associated vector space [3].

Word embeddings are calculated by ANNs trained on large corpora, and many sets of such embeddings for a lot of different languages exist. However, rather than using pre-calculated word embeddings, we found that including an Embedding layer in our models and calculate/learn from scratch the embeddings jointly

with the classification task produced better F1-scores on the dev. set.

In order to calculate the desired word embeddings, we first tokenize text, i.e. decompose it to individual words, symbols, punctuation marks etc. Each token is assigned an index and we consider a vocabulary of the most frequent tokens. Further, we set the length of the text’s representation as a sequence of tokens to a fixed length. Both the vocabulary’s size and the text’s length are hyperparameters with which one can experiment.

### 2.2 Undersampling

As mentioned in [1] the dataset is skewed/imbalanced; there are fewer samples of the positive class (i.e. flood-related) than the negative (approximately 20% - 80%). This makes training the model hard because during training it is presented with more negative samples and consequently ‘learns’ better the negative class and misclassifies a lot of positive samples, thus leading to a poor F1-score.

To tackle this issue, we use under sampling as follows: We keep all positive samples of the training set and select randomly some (not all) of the negative samples in order to have a set with a negative-positive class ratio closer to one and therefore a more balanced set. The value of this ratio is a hyperparameter which can be fine-tuned

### 2.3 ANN Models

Many ANN types for different tasks exist [2]. In this study, we are dealing with a binary classification problem whose solution may be viewed as a partition of the embeddings space into two sets, one for each class. This can be achieved by Multi-Layer Perceptron (MLP) added after the Embedding layer of the model. We chose a simple architecture of one hidden layer with 32 units having a ReLU activation function followed by a single output unit with a sigmoid activation function.

We then build on the previous model by considering a layer of the so-called Recurrent Neural Networks (RNN) consisting of 32 bidirectional LSTM units. RNNs are models where units have internal state acting as memory, thus they are capable of processing and learning sequence characteristics since they can ‘remember’ inputs seen in the past. A typical application of RNNs is time series prediction, but since text is a sequence of (correlated) words they are also used a lot in Natural Language Processing (NLP). The

LSTM layer is placed after the Embeddings layer and on top of that, we have the previous MLP structure.

Finally, we employed another type of ANN capable of handling sequences - the Convolutional Neural Network (CNN). Here learning a sequence is achieved via a different mechanism which exploits the mathematical operation of convolution of the input sequence with a small kernel. We thus placed after the Embeddings layer two parallel layers with 32 kernels of length 5 each. The outputs of those parallel Convolutional layers are then merged and being fed into the previous MLP architecture.

To convert the continuous (between zero and one) ANN output to binary (i.e. flood-related input text or not) we use a threshold. Texts having output above the threshold are labelled as flood-related (i.e. one) and texts having output below the threshold as labelled zero. The threshold is chosen for each model separately by maximizing the F1-score. Finally, the text's class was assigned by a majority rule on the three models' output.

### 3 RESULTS AND DISCUSSION

#### 3.1 Model setup and performance

After experimenting with various values, we ended up with a vocabulary of size 3000, sequence length of 40, embedding vector dimension of 300 and under-sampling ratio of 1.75. The vocabulary size and sequence length are small compared to typical Natural Language Processing (NLP) applications due to the short form of the tweet's text. The architecture of the ANNs used is described above.

ANNs were trained and evaluated individually on the same train/validation sets which were created by splitting the devset to an 80-20% ratio. The F1-scores on the validation set were 0.59 for the MLP, 0.60 for the RNN and CNN. Those scores were obtained by choosing thresholds 0.40, 0.65, 0.40 respectively. Finally, we combined the three ANN outputs by assigning to each input the majority class for the three ANN outputs. We chose this strategy, hoping that each ANN would perhaps capture different idiosyncrasies of the input. The overall F1 score improved slightly to 0.61. Our score on the test set was 0.5405, significantly lower, suggesting that we overfitted the training set.

#### 3.2 Limitations of the study

The main challenge of the task was related to the labelling of the training dataset. We noticed that many samples looked flood-related from a visual inspection but were not labeled as such (some example ids are:940319294084202496, 944240672294531073, 950753737466830940, 1059017654088790018, 1055172135587536896). Further, we noticed that many positive samples are from meteorological alerts. This could maybe restrict the training set and explain the difficulties of the model in generalizing well and thus, influence the overall model performance.

#### 3.3 Outlook - Ways to improve the performance

Experimenting with simpler text representations such as Bag of Words (BOW) and Term Frequency Inverse Document Frequency (TF-IDF) vectors and a Logistic Regression classifier revealed that taking into account tweet entities such as hashtags, in addition to the plain text, improved predictive performance.

However, due to time limitations, this approach was not implemented in our ANN framework. Further, it would require more sophisticated tokenization schemes able to extract hashtags, than those used for the ANNs input.

Geographical information of tweets, either in the form of metadata (e.g. coordinates, place attribute) or location mentions in the tweet's text could be exploited to 'geo locate' the tweet and possibly be used as additional inputs to the model. Especially since the dev. set focuses on a particular study area [1].

Finally, let us mention that this study focused solely on the tweet's text without considering the associated image. A two-branch model, where one branch would be the model presented here excluding the output layer and the other branch an image classifier both feeding the same output layer could be used to handle both text and image input.

#### 3.4 Code availability

The model was implemented as a Google Colab Ipython notebook and code is available upon request (theo\_nikolettopoulos@yahoo.co.uk).

### REFERENCES

- [1] Stelios Andreadis, Ilias Gialampoukidis, Anastasios Karakostas, Stefanos Vrochidis, Ioannis Kompatsiaris, Roberto Fiorin, Daniele Norbiato, and Michele Ferri. 2020. The Flood-related Multimedia Task at MediaEval 2020. In MediaEval 2020.
- [2] Ian Goodfellow, Yoshua Bengio, and Aaron Courville: Deep learning. [www.deeplearningbook.org](http://www.deeplearningbook.org)
- [3] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S. & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *Advances in neural information processing systems* (p./pp. 3111--3119)