

Gated Recurrent Neural Network Approach for Multilabel Emotion Detection in Microblogs

Prabod Rathnayaka¹, Supun Abeysinghe¹, Chamod Samarajeewa¹,
Isura Manchanayake¹, Malaka J. Walpola¹, Rashmika Nawaratne²,
Tharindu Bandaragoda² and Damminda Alahakoon²

¹Department of Computer Science and Engineering,
University of Moratuwa, Sri Lanka.

²Research Centre for Data Analytics and Cognition,
La Trobe University Victoria, Australia.

¹{prabod.14, supun.14, chamod.14, isura.14, malaka}@cse.mrt.ac.lk
²{b.nawaratne, t.bandaragoda, d.alahakoon}@latrobe.edu.au

Abstract

People express their opinions and emotions freely in social media posts and online reviews that contain valuable feedback for multiple stakeholders such as businesses and political campaigns. Manually extracting opinions and emotions from large volumes of such posts is an impossible task. Therefore, automated processing of these posts to extract opinions and emotions is an important research problem. However, human emotion detection is a challenging task due to the complexity and nuanced nature. To overcome these barriers, researchers have extensively used techniques such as deep learning, distant supervision, and transfer learning. In this paper, we propose a novel Pyramid Attention Network (PAN) based model for emotion detection in microblogs. The main advantage of our approach is that PAN has the capability to evaluate sentences in different perspectives to capture multiple emotions existing in a single text. The proposed model was evaluated on a recently released dataset and the results achieved the state-of-the-art accuracy of 58.9%.

1 Introduction

Emotions are an integral part of human life which in turn affects human decision-making process and human thinking patterns. People often tend to express their opinions freely in social media compared to other means. Typically microblogs written by social media users tend to have effects from their emotions towards the topic they discuss or the emotions they feel at the moment. Hence, the linguistic features of these texts are highly dependent on the emotion and therefore can be used to extract the underlying emotion. Identifying underlying emotions of microblogs is useful in understanding author's opinions. This becomes beneficial in natural language applications in diverse

fields including marketing, political campaigns, governing and human behavioral analysis.

Sentiment analysis can be considered as a fundamental type of emotion analysis. Though there have been a significant volume of research about sentiment analysis over the years, research about emotion detection and analysis has not gained much attention. A potential reason for that is complex and subtle behavior of human emotions compared to simple negative/positive sentiments. Research literature shows that there have been attempts to tackle this challenge using distant supervision techniques where emojis or hashtags present in the text are considered as indicators of emotions. However, such data can be noisy and somewhat unreliable, which affects the accuracy of such approaches.

Recently [Mohammad et al. \(2018\)](#) has released a significantly large dataset as a SemEval 2018 task. Recent advancements in deep learning for natural language processing shows, given enough data, it is often possible to develop a model which achieves a reasonable level of accuracy. Frequently, attention mechanisms are employed in such deep learning models. In this paper, we propose a novel attention mechanism *Pyramid Attention Network*(PAN) which has the ability to attend sentences from different perspectives. This becomes vital in emotion detection since a single micro-blog can contain multiple emotions and it needs to be considered in different perspectives to extract all the inherent emotions. State-of-the-art results achieved in our experiments is a good indication of the effectiveness of this approach.

2 Related Work

[Mohammad \(2012\)](#), [Mohammad and Kiritchenko \(2015\)](#), [Wang et al. \(2012\)](#), [Volkova and Bachrach \(2016\)](#) and [Abdul-Mageed and Ungar \(2017\)](#) used

distant supervision learning approach to acquire Twitter data with emotion related hashtags and emoticons. Abdul-Mageed and Ungar (2017) presents the largest dataset out of above all and they conducted a validation of the collected dataset using human annotators. A sample of 5600 tweets from the dataset has been checked and the study revealed only 61.37% of the inferred emotions were relevant. Hence model evaluations for emotion detection done using such noisy datasets can be inaccurate.

Another approach for emotion detection is to use transfer learning. The frequent usage of emojis can be seen in social media posts to express the emotion associated with the text. DeepMoji (Felbo et al., 2017) has addressed emotion prediction using a proposed variant of transfer learning called ‘chain-thaw’ where they used a pre-trained model which predicts emoji occurrences in microblogs. The model is based on a bidirectional Long Short Term Memory (LSTM) Network (Hochreiter and Schmidhuber, 1997) which was trained using a huge dataset of 1 billion tweets.

LSTMs (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Units (GRU) (Cho et al., 2014; Chung et al., 2015) which are variants of Recurrent Neural Networks (RNN), and Convolutional Neural Networks (CNN) have shown state-of-the-art results in text classification tasks including sentiment analysis (Ren et al., 2016; Liu et al., 2015; Tai et al., 2015; Tang et al., 2015; Zhang et al., 2016; Kalchbrenner et al., 2014; Kim, 2014; Zhang et al., 2015). Baziotis et al. (2018) got top results in the SemEval task using transfer learning approach combined with Deep Attentive RNNs (Bahdanau et al., 2014). Among other top results, Park et al. (2018) used a transfer learning based approach whereas Kim et al. (2018) and Rozen-tal and Fleischer (2018) used attention based approaches.

3 Methodology

First, the raw tweets are preprocessed, then mapped into a continuous vector space using an embedding layer. Then two stacked layers of Bidirectional GRUs (Cho et al., 2014) are used for feature extraction. This is followed by the novel attention mechanism where weighted average of the extracted features are taken. Finally, a dense layer maps this to eleven emotion categories using a sig-

moid activation function.

3.1 Preprocessing

Preprocessing is a key step in the model as it affects the accuracy of the model significantly. We have used ekphrasis tool introduced by Baziotis et al. (2017) for tweet preprocessing. Tweet tokenizing, word normalization, spell correcting and word segmentation for hashtags are done as pre-processing steps.

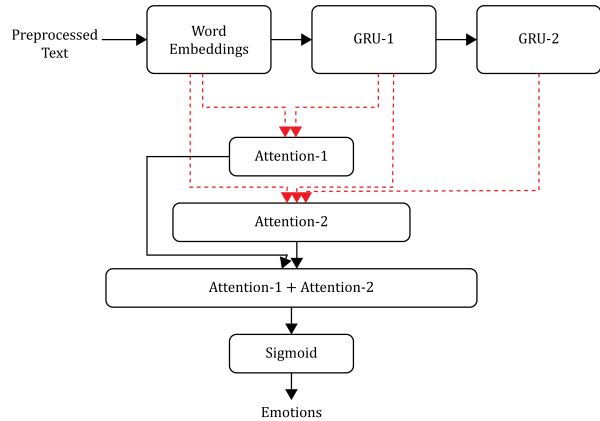


Figure 1: Overall model architecture

3.2 Embedding Layer

Given a sentence $S = [s_1, s_2, \dots, s_t, \dots, s_T]$ where s_t is the one hot vector representation of word at position t , the words are embedded into a continuous vector space using an embedding matrix W_e .

$$X = SW_e \quad (1)$$

X serves as the embedding of the input sequence and is fed to the first GRU Layer. We used pre-trained Glove (Pennington et al., 2014) word embeddings (300 dimension) in our model. These pre-trained word vectors are kept frozen so that they are not updated during back-propagation. Using pre-trained word embeddings improved the model performance notably.

3.3 Bidirectional GRU Layers

Gated Recurrent Unit (GRU) (Cho et al., 2014) is an improved version of a standard Recurrent Neural Network designed to overcome the vanishing gradient problem and exploding problem (Bengio et al., 1994; Hochreiter and Schmidhuber, 1997). GRU uses two special gating mechanisms called, update gate (z_t) and reset gate (r_t) to decide which information is passed to the output. Update gate

helps the model to determine how much of past information is passed on to the future steps. Reset gate helps the model to determine how much of past information to forget. Then the output of position t is created using concatenating forward and backward hidden states ($\vec{h}_t, \overleftarrow{h}_t$). Hidden size (output size) of both GRU layers is set to 50. W and b in (2), (3) and (4) corresponds to the weights and biases for the gates.

$$r_t = \sigma(W_{ir}x_t + b_{ir} + W_{hr}h_{(t-1)} + b_{hr}) \quad (2)$$

$$z_t = \sigma(W_{iz}x_t + b_{iz} + W_{hz}h_{(t-1)} + b_{hz}) \quad (3)$$

$$n_t = \tanh(W_{in}x_t + b_{in} + r_t(W_{hn}h_{(t-1)} + b_{hn})) \quad (4)$$

$$h_t = (1 - z_t)n_t + z_th_{(t-1)} \quad (5)$$

3.4 Pyramid Attention Network

Neural Networks with Attention has shown success in a variety of tasks such as machine translation (Luong et al., 2015), question answering (Yang et al., 2016a), image captioning (You et al., 2016) and sentiment analysis (Yang et al., 2016b). Aforementioned tasks such as sentiment analysis which can infer the sentiment based on a set of keywords, can directly use attention to focus on a specific part of a sentence and decide the sentiment. However, in multi-label emotion classification, attending to a part of a sentence can infer some emotions but could not capture every emotion embedded in the sentence. Following is a tweet extracted from the dataset.

The best revenge is massive success. →anger, joy, optimism

By giving high attention to “revenge”, we can infer the emotion *anger*, but to infer emotions *joy* and *optimism* we need to consider the whole sentence in a different perspective. Therefore, a single attention vector will not give us all the emotions associated.

To overcome the above limitation, we propose a new attention architecture, which we call *Pyramid Attention Network*. In the proposed model, attention layer-1 attends to the word embeddings and GRU layer-1 whereas attention layer-2 attends to both of them and GRU layer-2 in addition. This can be generalized to a n-layer stacked GRU or LSTM as well.

The fundamental concept behind this attention mechanism is inspired by Bahdanau et al. (2014), Yang et al. (2016b) and Felbo et al. (2017) where the key difference lies in the attention architecture

in general. Felbo et al. (2017) follows a similar approach but only limits to a single attention layer which attends to all the previous layers. Yang et al. (2016b) used two attention layers but it contrasts from our model since it uses a word attention layer followed by a sentence embedding layer and subsequently a sentence attention layer. Our model is formally defined as follows.

For an arbitrary vector U we can define attention coefficients (a_i) and output V as follows.

Let $u_i \in U$,

$$V = \sum_i a_i u_i \quad (6)$$

$$\text{where, } e_i = u_i w_a + b \text{ and } a_i = \frac{\exp(e_i)}{\sum_t \exp(e_t)} \quad (7)$$

For the first attention layer, $V_1 = (H_1, X)$, and for the second attention layer, $V_2 = (H_2, H_1, X)$ where H_1 and H_2 are outputs of Bi-GRU layer 1 and layer 2 respectively and X is the word embeddings for the input sentence.

3.5 Classification and Training

Output of the attention layers 1 (V_1) and 2 (V_2) are concatenated ($V = (V_1, V_2)$) and passed into a dense layer with a sigmoid activation which outputs a vector of size 11. It has a value between 0 and 1 for each emotion class. If the value is larger than a threshold value, it is classified as positive. We used 0.5 as this threshold.

$$\hat{y} = \sigma(W_d V + b) \quad (8)$$

Weighted binary cross entropy loss function is used with a weight of $w = 2$ for the correctly labeled ones. y represents the ground truth labels and \hat{y} represents the predicted values. m is the number of emotions, which is 11 in this scenario.

$$J(\theta) = -\frac{1}{m} \sum_{i=0}^m (wy_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)) \quad (9)$$

Regularization is essential to ensure that model is not over-fitted in the training phase. Dropout is a regularization method proposed by Srivastava et al. (2014) which randomly turn-off a percentage of the neurons of a layer in the network. We apply dropout of 0.2 between the attention layer and the dense layer, spatial dropout of 0.4 between word embedding layer and GRU layer for regularization. We added a Gaussian Noise with 0.1 standard deviation to GRU hidden weights to reduce

over-fitting. In addition, we apply L2 regularization penalty to the loss function to reduce large weights. Moreover, we used early stopping (Caruana et al., 2001) where training is stopped after the validation loss stops decreasing. The model was trained to minimize the weighted binary cross entropy loss using backpropagation. We used Adam optimizer (Kingma and Ba, 2014) with a batch size of 64 and an initial learning rate of 0.001, which will be reduced by half for every 3 consecutive failures to reduce validation loss with a lower bound of 0.0001. All the values for regularizations were found empirically.

4 Experiments and Results

We used the recently published SemEval 2018 Task 1 emotion classification dataset (Mohammad et al., 2018) for the evaluation of the model. This dataset consists of 10983 tweets which are categorized into 11 emotion categories; *anger*, *anticipation*, *disgust*, *fear*, *joy*, *love*, *optimism*, *pessimism*, *sadness*, *surprise* and *trust*. Each tweet can have multiple emotions, thus multi-label classification mechanisms are employed. As per the SemEval task, the dataset is split into training, development, and testing set with respectively 6838, 886 and 3259 tweets for each set.

We have compared our proposed model with the top four models of the SemEval 2018 Task 1 (Baziotis et al., 2018; Park et al., 2018; Kim et al., 2018; Rozental and Fleischer, 2018) and our implementation of Hierarchical Attention Network (HAN) approach proposed by Yang et al. (2016b) adopted to multi-label case. As shown in Table 1, proposed model achieves the state-of-the-art for the emotion detection dataset by Mohammad et al. (2018).

5 Analysis of Results

Figure 2 shows the number of tweets contained in each category vs the measured F-score achieved by our proposed model for each category. Six categories stand out the most with a significantly large amount of data whereas rest of the categories contain a comparably smaller number of tweets. It further shows tweets with a significant number of training examples are detected with a higher F-score whereas underrepresented categories have a lower F-score. Such underrepresented emotions affect the performance of the proposed model severely. For other well represented

Model	Jaccard	Micro	Macro
Baziotis et al. (2018)	0.579	-	-
Park et al. (2018)	0.576	0.692	0.497
Kim et al. (2018)	0.574	0.687	0.511
Rozental and Fleischer (2018)	0.566	0.673	0.490
HAN (Yang et al., 2016b)	0.567	0.683	0.535
Proposed Model	0.589	0.701	0.550

Table 1: Comparing results of the proposed model. Proposed model achieves state-of-the-art for the dataset by Mohammad et al. (2018). **Jaccard** - mutli-label accuracy (Jaccard accuracy), **Micro** - Micro-avg F1 score and **Macro** - Macro-avg F1 score.

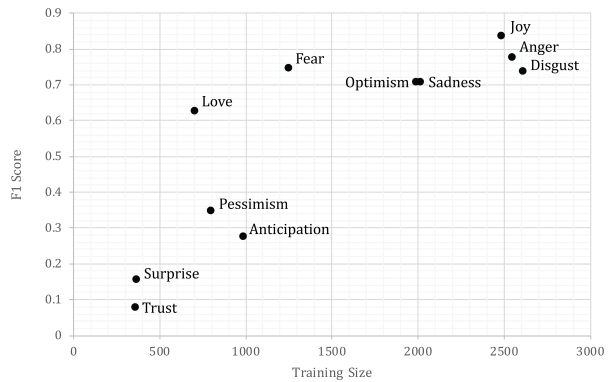


Figure 2: The variation of F1 score with the dataset size of individual emotions

emotions, our model achieves a F-score more than 0.7. One interesting fact to note is that, though there are comparably smaller number of tweets containing *love*, still our proposed model managed to predict that emotion accurately. The reason for this can be *love* is highly correlated with special keywords.

6 Conclusion and Future Work

We proposed a novel *Pyramid Attention Network* (PAN) which achieves state-of-the-art performance for emotion detection in microblogs. Our analysis revealed that there are underrepresented categories in the dataset and those classes affect the model accuracy significantly. Future work includes enhancing the overall performance by improving the detection of underrepresented classes, experiment this approach for multiple emotion detection datasets, and apply this attention mechanism to other text classification tasks.

References

- Muhammad Abdul-Mageed and Lyle Ungar. 2017. Emonet: Fine-grained emotion detection with gated recurrent neural networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 718–728.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Christos Baziotis, Nikos Athanasiou, Alexandra Chronopoulou, Athanasia Kolovou, Georgios Paraskevopoulos, Nikolaos Ellinas, Shrikanth Narayanan, and Alexandros Potamianos. 2018. Ntua-slp at semeval-2018 task 1: Predicting affective content in tweets with deep attentive rnns and transfer learning. *arXiv preprint arXiv:1804.06658*.
- Christos Baziotis, Nikos Pelekis, and Christos Douk-eridis. 2017. Datastories at semeval-2017 task 4: Deep lstm with attention for message-level and topic-based sentiment analysis. In *Proceedings of the 11th International Workshop on Semantic Evaluation (SemEval-2017)*, pages 747–754, Vancouver, Canada. Association for Computational Linguistics.
- Yoshua Bengio, Patrice Simard, and Paolo Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks*, 5(2):157–166.
- Rich Caruana, Steve Lawrence, and C Lee Giles. 2001. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *Advances in neural information processing systems*, pages 402–408.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*.
- Junyoung Chung, Caglar Gulcehre, Kyunghyun Cho, and Yoshua Bengio. 2015. Gated feedback recurrent neural networks. In *International Conference on Machine Learning*, pages 2067–2075.
- Bjarke Felbo, Alan Mislove, Anders Søgaard, Iyad Rahwan, and Sune Lehmann. 2017. Using millions of emoji occurrences to learn any-domain representations for detecting sentiment, emotion and sarcasm. *arXiv preprint arXiv:1708.00524*.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. 2014. A convolutional neural network for modelling sentences. *arXiv preprint arXiv:1404.2188*.
- Yanghoon Kim, Hwanhee Lee, and Kyomin Jung. 2018. Attnconvnet at semeval-2018 task 1: Attention-based convolutional neural networks for multi-label emotion classification. *arXiv preprint arXiv:1804.00831*.
- Yoon Kim. 2014. Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Pengfei Liu, Xipeng Qiu, Xinchu Chen, Shiyu Wu, and Xuanjing Huang. 2015. Multi-timescale long short-term memory neural network for modelling sentences and documents. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 2326–2335.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. 2015. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*.
- Saif M Mohammad. 2012. # emotional tweets. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*, pages 246–255. Association for Computational Linguistics.
- Saif M Mohammad, Felipe Bravo-Marquez, Mohammad Salameh, and Svetlana Kiritchenko. 2018. Semeval-2018 task 1: Affect in tweets. In *Proceedings of International Workshop on Semantic Evaluation (SemEval-2018)*, New Orleans, LA, USA.
- Saif M Mohammad and Svetlana Kiritchenko. 2015. Using hashtags to capture fine emotion categories from tweets. *Computational Intelligence*, 31(2):301–326.
- Ji Ho Park, Peng Xu, and Pascale Fung. 2018. Plusemo2vec at semeval-2018 task 1: Exploiting emotion knowledge from emoji and# hashtags. *arXiv preprint arXiv:1804.08280*.
- Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. Glove: Global vectors for word representation. In *Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543.
- Yafeng Ren, Yue Zhang, Meishan Zhang, and Donghong Ji. 2016. Context-sensitive twitter sentiment classification using neural network. In *AAAI*, pages 215–221.
- Alon Rozental and Daniel Fleischer. 2018. Amobee at semeval-2018 task 1: Gru neural network with a cnn attention mechanism for sentiment classification. *arXiv preprint arXiv:1804.04380*.

- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958.
- Kai Sheng Tai, Richard Socher, and Christopher D Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. *arXiv preprint arXiv:1503.00075*.
- Duyu Tang, Bing Qin, and Ting Liu. 2015. Document modeling with gated recurrent neural network for sentiment classification. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 1422–1432.
- Svitlana Volkova and Yoram Bachrach. 2016. Inferring perceived demographics from user emotional tone and user-environment emotional contrast. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1567–1578.
- Wenbo Wang, Lu Chen, Krishnaprasad Thirunarayan, and Amit P Sheth. 2012. Harnessing twitter” big data” for automatic emotion identification. In *Privacy, Security, Risk and Trust (PASSAT), 2012 International Conference on and 2012 International Conference on Social Computing (Social-Com)*, pages 587–592. IEEE.
- Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. 2016a. Stacked attention networks for image question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 21–29.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016b. Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1480–1489.
- Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. 2016. Image captioning with semantic attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4651–4659.
- Meishan Zhang, Yue Zhang, and Duy-Tin Vo. 2016. Gated neural networks for targeted sentiment analysis. In *AAAI*, pages 3087–3093.
- Xiang Zhang, Junbo Zhao, and Yann LeCun. 2015. Character-level convolutional networks for text classification. In *Advances in neural information processing systems*, pages 649–657.