

SAME SIDE STANCE CLASSIFICATION TASK: Facilitating Argument Stance Classification by Fine-tuning a BERT Model

Stefan Ollinger, Lorik Dumani, Premtim Sahitaj, Ralph Bergmann, Ralf Schenkel

University of Trier

D-54286 Trier

stefan.ollinger@gmx.de

{dumani, s4prsa, bergmann, schenkel}@uni-trier.de

Abstract

Research on computational argumentation is currently being intensively investigated. The goal of this community is to find the best pro and con arguments for a user given topic either to form an opinion for oneself, or to persuade others to adopt a certain standpoint. While existing argument mining methods can find appropriate arguments for a topic, a correct classification into pro and con is not yet reliable. The same side stance classification task provides a dataset of argument pairs classified by whether or not both arguments share the same stance and does not need to distinguish between topic-specific pro and con vocabulary but only the argument similarity within a stance needs to be assessed. The results of our contribution to the task are build on a setup based on the BERT architecture. We fine-tuned a pre-trained BERT model for three epochs and used the first 512 tokens of each argument to predict if two arguments share the same stance.

1 Introduction

Argumentation is an activity in everyday human life. We argue in domains such as health, law and politics either trying to find standpoints which are acceptable by being supported by reasons or to persuade others to a certain point of view and if necessary to carry out certain actions. Computational Argumentation (CA) aims to find argument representations and models which are well suited to do computation with arguments. CA is a new and fast growing field of research. In the simplest case, an *argument* is a *claim* supported or opposed by at least one *premise* (Peldszus and Stede, 2013). An example of a claim c could be “We need to abolish nuclear energy”, examples of premises that support and oppose this claim could be $p_1 =$ “Renewable energy sources will eventually be able to replace fossil fuel and nuclear en-

ergy” and $p_2 =$ “Nuclear energy is a cheap alternative to fossil fuels”, respectively. Common tasks in CA include argument mining (AM) and argument retrieval (AR). AM reconstructs arguments from textual sources, e.g. in form of an argument graph. AR finds all relevant arguments for a topic. Existing argument search engines like ARGS¹ or ARGUMENTEXT² search for the best supporting and opposing premises for a user query on a usually controversial topic and list them separately in pro and con. The correct classification of stances is therefore a fundamental task in computational argumentation. However, a short-coming of current stance classification algorithms is that their classifiers must be trained for a particular topic, i.e. they cannot be reliably applied across topics (Webis, 2019a).

In the task \gg SAME SIDE (STANCE) CLASSIFICATION \ll a simplified variant is to be examined, namely whether two given arguments to a topic have the same stance. For example, p_1 and p_2 have different stances, but p_1 and $p_3 =$ “The danger from radioactive contamination should be avoided” would have the same stance to the topic *nuclear energy*. The particular difficulty lies in the fact that p_1 and p_3 are syntactically very different. So we have to decide on a semantic level whether the stances are the same.

In this paper we present a method where we fine-tune a pre-trained BERT (Devlin et al., 2019) model to decide whether two arguments have the same stance. In Section 2 we discuss related work. In Section 3 we specify the dataset. Then, in Section 4 we describe the implementation and the evaluation of our approach. Finally, Section 5 concludes the paper.

¹www.args.me

²www.argumenttext.de

2 Related Work

In our implementation we make use of BERT (Devlin et al., 2019) which achieved state-of-the-art results in many NLP tasks such as Natural Language Inference (MNLI), semantic textual similarity (STS), and others. BERT makes use of the Transformer architecture (Vaswani et al., 2017), more precisely it applies a bidirectional masked language model training to the architecture. Contrary to previous embedding techniques like WORD2VEC (Mikolov et al., 2013) its mechanism learns contextualized representations of words in a text.

(Stab and Gurevych, 2014) address argumentative relation classification. The relation between two argument components is divided in *support* and *non-support* classes. Therefore a range of structural, lexical and syntactic features are defined and extracted for an argument component pair. The classification is done with a SVM.

Bar-Haim et al. (2017a) address stance classification of premises towards a claim topic. Here, the classification task is divided further into simpler sub-tasks. Bar-Haim et al. (2017b) extend this work by a more extensive sentiment lexicon and contextual features.

Most presented approaches classify the relation of a premise towards a claim. In contrast to the same stance classification the relation between two premises is considered. Further we do not apply feature engineering, but rely on the neural network to extract good features.

3 The provided Dataset

The arguments from the provided dataset were extracted from the four web sources idebate.org, debatepedia.org, debatewise.org and debate.org. Each instance consists of the seven fields that are depicted in Table 1.

The two most discussed topics “abortion” and “gay marriage” were chosen and two experiments were set-up for the same side stance classification task. The first experiment addresses the classification within topics and consists of a training set with arguments for a set of topics (abortion and gay marriage) and a test set with arguments related to the same set of topics. Table 2 illustrates an overview of the data within topics.

For the second experiment, which addresses the classification across topics, the training set contains arguments for a topic (abortion) and the test

Label	Description
<i>id</i>	The id of the instance
<i>topic</i>	The title of the debate. It can be a general topic (e.g. abortion) or a topic with a stance (e.g. abortion should be legalized).
<i>argument1</i>	A pro or con argument related to the topic.
<i>argument1_id</i>	The ID of argument1.
<i>argument2</i>	A pro or con argument related to the topic.
<i>argument2_id</i>	The ID of argument2.
<i>is_same_stance</i>	True or False. True in case argument1 and argument2 have the same stance towards the topic and False otherwise.

Table 1: Fields of each instance in the dataset. Source: (Webis, 2019a)

class	topic:	topic:
	abortion	gay marriage
<i>Same Side</i>	20,834	13,277
<i>Different Side</i>	20,006	9,786
<i>Total</i>	40,840	23,063

Table 2: Overview of the data within topics. Source: (Webis, 2019a)

set arguments are related to another set of topics. The class *Same Side* contains 31,195 instances, the class *Different Side* 29,853.

4 Evaluation

In this section we describe the experimental setup and evaluate our approach utilizing BERT and compare it to a SVM baseline.

4.1 Hypotheses

In order to measure the performance of our approach, the following hypotheses were formulated and are subject of this evaluation:

- **H1:** A Transformer-based sequence classification improves upon the SVM baseline.
- **H2:** The *large* Transformer model outperforms the smaller *base* model.
- **H3:** Longer input sequences yield better classification than shorter sequence lengths.

- **H4:** Classification of full sentences performs better than including partial input sentences.

4.2 Experimental Setup

First, we divided the provided data set into training and test sets (90% and 10%). This results in 57,512 training pairs and 6,391 test pairs for within topics classification as well as 54,943 training and 6,105 test pairs for cross topics taken all from the shared task labeled training data. Then we used BERT (Devlin et al., 2019)³ in our implementation for training to classify arguments of same stance. We used both provided models *base* and *large* always with three epochs for fine-tuning. All models use lower-cased token sequences and vocabulary. It should be noted here that BERT is limited to a fixed size of tokens, with the maximum being 512 tokens for the pre-trained models. Longer input sequences are truncated to the maximum sequence length. This truncation can lead to loss of information which we evaluate in hypotheses 3 and 4.

4.3 Results and Discussion

Figure 1 shows the results within the same topic with varying maximum sequence length. Figure 2 shows the results between topics. Argument pairs whose length exceeds the maximum sequence lengths are uniformly truncated. In within topic evaluation *large* yields higher accuracy than *base* in four of five cases. In the cross topic evaluation the result is similar. Therefore hypothesis 2 can be accepted. Nevertheless the smaller *base* model is quite close to *large*.

The SVM baseline, supplied by the shared task organizers, achieves 54% accuracy in within topic and 58% cross topic. The *base* model improves upon this already with the smallest sequence length of 32 tokens. Thus hypothesis 1 can be accepted. This result is possibly due to a Transformer having a larger model capacity and employing better suited representations for natural language text compared to an SVM.

Next, we take a look at argumentative input of varying maximum sequence length. We can see from Figure 1 and Figure 2 that the classification benefits from more contextual information. Hypothesis 3 can therefore be accepted. One question is why a model using 64 tokens already performs quite well. Figure 3 shows the distribution of the

³<https://github.com/huggingface/transformers>

summed argument lengths. We can observe here that the majority (76%) consists of less than 512 tokens. As we can infer from Figure 4, the distribution of the lengths is usually even considerably below 512 tokens. This explains why models with rather short contextual information perform quite well already.

Since the input sequences are truncated the Transformer model also trains with incomplete, partial natural language sentences. In order to see whether a model can better learn from full sentences we filter out all partial sequences which are longer than 512 tokens, reducing the available training/testing data (no trunc train/test). As can be seen in Table 3 and Table 4 the Transformer is able to learn from partial sentences. Therefore hypothesis 4 needs to be rejected. The highest results are achieved when testing is also done on untruncated full sentences. This result is an indicator of what could be achieved with Transformers of larger or variable maximum sequence length such as explored by (Dai et al., 2019).

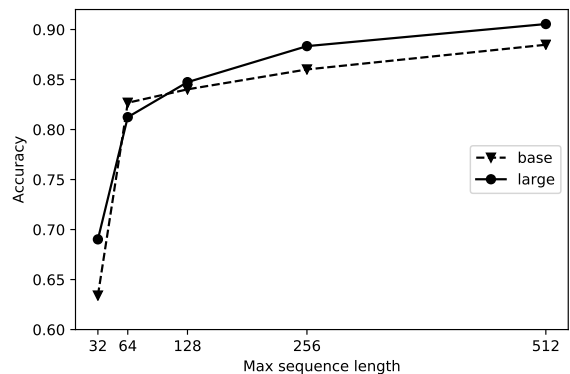


Figure 1: Accuracy of *base* and *large* model **within** topics for varying maximum sequence length

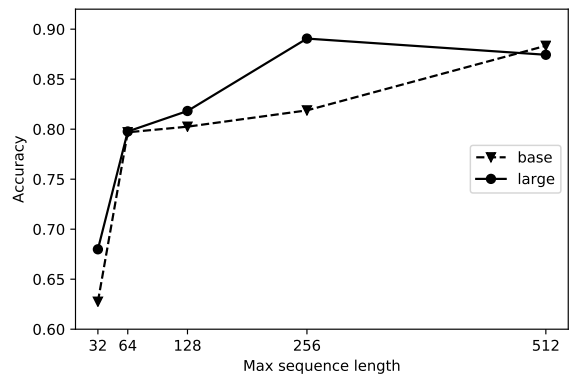


Figure 2: Accuracy of *base* and *large* model **across** topics for varying maximum sequence length

Table 3: Results of truncated/non-truncated training/testing within topics.

No Trunc Train	No Trunc Test	Model	# Train	# Test	Acc	F1
		bert-base	57,512	6,391	0.8848	0.8911
		bert-large	57,512	6,391	0.9055	0.9104
•		bert-base	43,678	6,391	0.8603	0.8668
•		bert-large	43,678	6,391	0.8830	0.8896
•	•	bert-base	43,678	4,932	0.9064	0.9137
•	•	bert-large	43,678	4,932	0.9471	0.9527

Table 4: Results of truncated/non-truncated training/testing across topics.

No Trunc Train	No Trunc Test	Model	# Train	# Test	Acc	F1
		bert-base	54,943	6,105	0.8834	0.8848
		bert-large	54,943	6,105	0.8744	0.8757
•		bert-base	40,763	6,105	0.8139	0.7971
•		bert-large	40,763	6,105	0.8622	0.8667
•	•	bert-base	40,763	4,515	0.8997	0.9026
•	•	bert-large	40,763	4,515	0.9271	0.9325

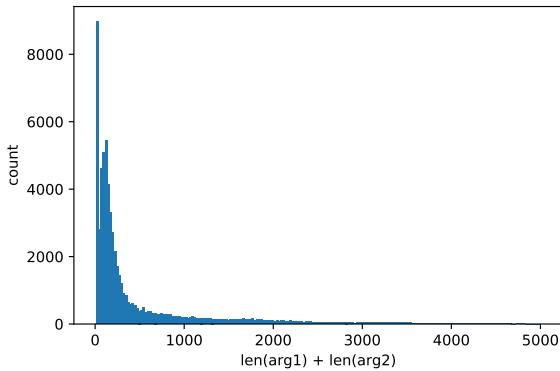


Figure 3: Distribution of argument pair lengths

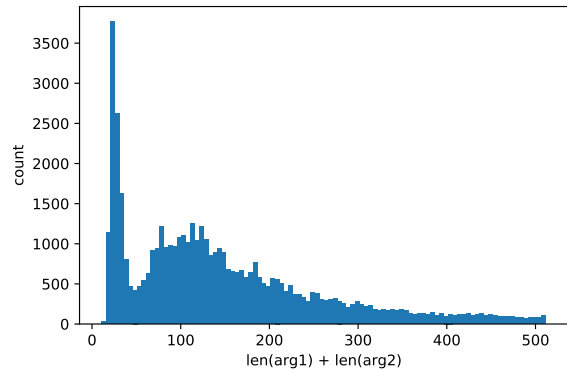


Figure 4: Distribution of token lengths within a maximum of 512 tokens.

5 Conclusion

In this paper we have contributed to the \gg SAME SIDE (STANCE) CLASSIFICATION \ll task and proposed a method which uses a fine-tuned BERT model to determine whether two given arguments have the same stance. The baseline of the organizers was outperformed with our method. In our evaluation the large model performs better than the base model. Our results also show that longer input sentences are classified better than shorter ones, and that classifying whole sentences does not perform better than classifying partial sentences. According to the organizers’ leaderboard (Webis, 2019b)⁴ our approach performed

best across topics with precision and recall values of 0.72 and an accuracy of 0.73. For within-topics we achieved the best performance as well as the ASV team from Leipzig University with an accuracy of 0.77. However, for this task we had a higher precision (0.85 vs. 0.79) but a lower recall (0.66 vs. 0.73).

Acknowledgements

This work has been funded by the Deutsche Forschungsgemeinschaft (DFG) within the project ReCAP, Grant Number 375342983 - 2018-2020, as part of the Priority Program “Robust Argumentation Machines (RATIO)” (SPP-1999).

⁴Ranking on the 16th August 2019.

References

- Roy Bar-Haim, Indrajit Bhattacharya, Francesco Dinuzzo, Amrita Saha, and Noam Slonim. 2017a. [Stance classification of context-dependent claims](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017*, pages 251–261.
- Roy Bar-Haim, Lilach Edelstein, Charles Jochim, and Noam Slonim. 2017b. [Improving claim stance classification with lexical knowledge expansion and context utilization](#). In *Proceedings of the 4th Workshop on Argument Mining, ArgMining@EMNLP 2017*, pages 32–38.
- Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc Viet Le, and Ruslan Salakhutdinov. 2019. [Transformer-xl: Attentive language models beyond a fixed-length context](#). In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019*, pages 2978–2988.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Gregory S. Corrado, and Jeffrey Dean. 2013. [Distributed representations of words and phrases and their compositionality](#). In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013*, pages 3111–3119.
- Andreas Peldszus and Manfred Stede. 2013. [From argument diagrams to argumentation mining in texts: A survey](#). *International Journal of Cognitive Informatics and Natural Intelligence*, 7(1):1–31.
- Christian Stab and Iryna Gurevych. 2014. [Identifying argumentative discourse structures in persuasive essays](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 46–56.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 5998–6008.
- Webis. 2019a. Same Side Stance Classification. <https://sameside.webis.de/>. Accessed: 2019-12-06.
- Webis. 2019b. Same Side Stance Classification Leaderboard. <https://sameside.webis.de/leaderboard.html>. Accessed: 2019-12-06.