
Experiments in Detecting Persuasion Techniques in the News

Seunghak Yu
MIT CSAIL
Cambridge, MA, USA
seunghak@csail.mit.edu

Giovanni Da San Martino
Qatar Computing Research Institute, HBKU
Doha, Qatar
gmartino@hbku.edu.qa

Preslav Nakov
Qatar Computing Research Institute, HBKU
Doha, Qatar
pnakov@hbku.edu.qa

Abstract

Many recent political events, like the 2016 US Presidential elections or the 2018 Brazilian elections have raised the attention of institutions and of the general public on the role of Internet and social media in influencing the outcome of these events. We argue that a safe democracy is one in which citizens have tools to make them aware of propaganda campaigns. We propose a novel task: performing fine-grained analysis of texts by detecting all fragments that contain propaganda techniques as well as their type. We further design a novel multi-granularity neural network, and we show that it outperforms several strong BERT-based baselines.

1 Introduction

Journalistic organisations, such as *Media Bias/Fact Check*,¹ provide reports on news sources highlighting the ones that are propagandistic. Obviously, such analysis is time-consuming and possibly biased and it cannot be applied to the enormous amount of news that flood social media and the Internet. Research on detecting propaganda has focused primarily on classifying entire articles as propagandistic/non-propagandistic [1, 2, 11]. Such learning systems are trained using gold labels obtained by transferring the label of the media source, as per *Media Bias/Fact Check* judgment, to each of its articles. Such distant supervision setting inevitably introduces noise in the learning process [8] and the resulting systems tend to lack explainability.

We argue that in order to study propaganda in a sound and reliable way, we need to rely on high-quality trusted professional annotations and it is best to do so at the fragment level, targeting specific techniques rather than using a label for an entire document or an entire news outlet. Therefore, we propose a novel task: identifying specific instances of propaganda techniques used within an article. In particular, we design a novel multi-granularity neural network, and we show that it outperforms several strong BERT-based baselines.

Our corpus could enable research in propagandistic and non-objective news, including the development of explainable AI systems. A system that can detect instances of use of specific propagandistic techniques would be able to make it explicit to the users why a given article was predicted to be propagandistic. It could also help train the users to spot the use of such techniques in the news.

¹<http://mediabiasfactcheck.com/>

2 Corpus Annotated with Propaganda Techniques

We retrieved 451 news articles from 48 news outlets, both propagandistic and non-propagandistic according to *Media Bias/Fact Check*, which professionals annotators² annotated according to eighteen persuasion techniques [9], ranging from leveraging on the emotions of the audience —such as using *loaded language* or *appeal to authority* [6] and slogans [4]— to using logical fallacies —such as *straw man* [12] (misrepresenting someone’s opinion), hidden *ad-hominem fallacies*, and *red herring* [13, p. 78] (presenting irrelevant data).³ Some of these techniques weren studied in tasks such as hate speech detection and computational argumentation [7].

The total number of technique instances found in the articles, after the consolidation phase, is 7,485, out of a total number of 21,230 sentences (35.2%). The distribution of the techniques in the corpus is also uneven: while there are 2,547 occurrences of *loaded language*, there are only 15 instances of *straw man* (more statistics about the corpus can be found in [3]). We define two tasks based on the corpus described in Section 2: (i) **SLC (Sentence-level Classification)**, which asks to predict whether a sentence contains at least one propaganda technique, and (ii) **FLC (Fragment-level classification)**, which asks to identify both the spans and the type of propaganda technique. Note that these two tasks are of different granularity, g_1 and g_2 , namely tokens for FLC and sentences for SLC. We split the corpus into training, development and test, each containing 293, 57, 101 articles and 14,857, 2,108, 4,265 sentences, respectively.

Our task requires specific evaluation measures that give credit for partial overlaps of fragments. Thus, in our precision and recall versions, we give partial credit to imperfect matches at the character level, as in plagiarism detection [10].

Let s and t be two fragments, i.e., sequences of characters. We measure the overlap of two annotated fragments as $C(s, t, h) = \frac{|(s \cap t)|}{h} \delta(l(s), l(t))$, where h is a normalizing factor, $l(a)$ is the labelling of fragment a , and $\delta(a, b) = 1$ if $a = b$, and 0 otherwise.

We now define variants of precision and recall able to account for the imbalance in the corpus:

$$P(S, T) = \frac{1}{|S|} \sum_{\substack{s \in S, \\ t \in T}} C(s, t, |s|), \quad R(S, T) = \frac{1}{|T|} \sum_{\substack{s \in S, \\ t \in T}} C(s, t, |t|), \quad (1)$$

In eq. (1), we define $P(S, T)$ to be zero if $|S| = 0$ and $R(S, T)$ to be zero if $|T| = 0$. Finally, we compute the harmonic mean of precision and recall in Eq. (1) and we obtain an F_1 -measure. Having a separate function C for comparing two annotations gives us additional flexibility compared to standard NER measures that operate at the token/character level, e.g., we can change the factor that gives credit for partial overlaps and be more forgiving when only a few characters are wrong.

3 Models

We depart from BERT [5], and we design three baselines.

BERT. We add a linear layer on top of BERT and we fine-tune it, as suggested in [5]. For the FLC task, we feed the final hidden representation for each token to a layer L_{g_2} that makes a 19-way classification: does this token belong to one of the eighteen propaganda techniques or to none of them (cf. Figure 1-a). For the SLC task, we feed the final hidden representation for the special [CLS] token, which BERT uses to represent the full sentence, to a two-dimensional layer L_{g_1} to make a binary classification.

BERT-Joint. We use the layers for both tasks in the BERT baseline, L_{g_1} and L_{g_2} , and we train for both FLC and SLC jointly (cf. Figure 1-b).

BERT-Granularity. We modify BERT-Joint to transfer information from SLC directly to FLC. Instead of using only the L_{g_2} layer for FLC, we concatenate L_{g_1} and L_{g_2} , and we add an extra 19-dimensional classification layer $L_{g_{1,2}}$ on top of that concatenation to perform the prediction for FLC (cf. Figure 1-c).

²<http://www.aiidatapro.com>. The company performs professional annotations in the NLP domain, although they were not expert in propaganda techniques before this work.

³For a complete list see <http://propaganda.qcri.org/annotations/definitions.html>

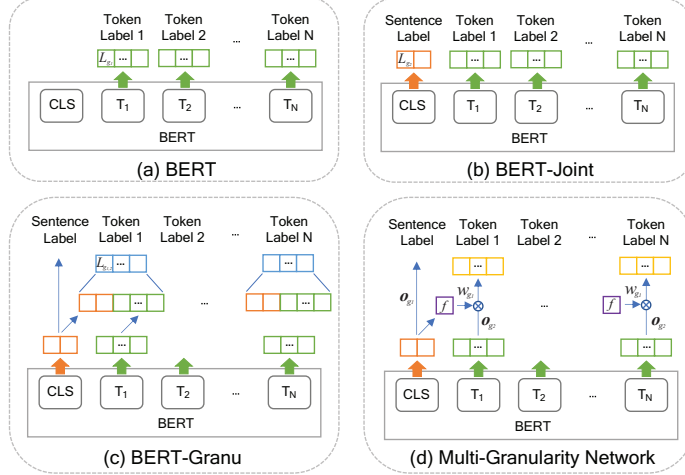


Figure 1: The architecture of the baseline models (a-c), and of our multi-granularity network (d).

Multi-Granularity Network. We propose a model that can drive the higher-granularity task (FLC) on the basis of the lower-granularity information (SLC), rather than simply using low-granularity information directly. Figure 1-d shows the architecture of this model.

More generally, suppose there are k tasks of increasing granularity, e.g., document-level, paragraph-level, sentence-level, word-level, subword-level, character-level. Each task has a separate classification layer L_{g_k} that receives the feature representation of the specific level of granularity g_k and outputs \mathbf{o}_{g_k} . The dimension of the representation depends on the embedding layer, while the dimension of the output depends on the number of classes in the task. The output \mathbf{o}_{g_k} is used to generate a weight for the next granularity task g_{k+1} through a trainable gate f :

$$w_{g_k} = f(\mathbf{o}_{g_k}) \quad (2)$$

The gate f consists of a projection layer to one dimension and an activation function. The resulting weight is multiplied by each element of the output of layer $L_{g_{k+1}}$ to produce the output for task g_{k+1} :

$$\mathbf{o}_{g_{k+1}} = w_{g_k} * \mathbf{o}_{g_{k+1}} \quad (3)$$

If $w_{g_k} = 0$ for a given example, the output of the next granularity task $\mathbf{o}_{g_{k+1}}$ would be 0 as well. In our setting, this means that, if the sentence-level classifier is confident that the sentence does not contain propaganda, i.e., $w_{g_k} = 0$, then $\mathbf{o}_{g_{k+1}} = 0$ and there would be no propagandistic technique predicted for any span within that sentence. Similarly, when back-propagating the error, if $w_{g_k} = 0$ for a given example, the final entropy loss would become zero, i.e., the model would not get any information from that example. As a result, only examples strongly classified as negative in a lower-granularity task would be ignored in the high-granularity task. Having the lower-granularity as the main task means that higher-granularity information can be selectively used as additional information to improve the performance, but only if the example is not considered as highly negative.

For the loss function, we use a cross-entropy loss with sigmoid activation for every layer, except for the highest-granularity layer L_{g_K} , which uses a cross-entropy loss with softmax activation. Unlike softmax, which normalizes over all dimensions, the sigmoid allows each output component of layer L_{g_k} to be independent from the rest. Thus, the output of the sigmoid for the positive class increases the degree of freedom by not affecting the negative class, and vice versa. As we have two tasks, we use sigmoid activation for L_{g_1} and softmax activation for L_{g_2} . Moreover, we use a weighted sum of losses with a hyper-parameter α :

$$\mathcal{L}_{\mathcal{J}} = \mathcal{L}_{g_1} * \alpha + \mathcal{L}_{g_2} * (1 - \alpha) \quad (4)$$

Again, we use BERT [5] for the contextualized embedding layer and we place the multi-granularity network on top of it.

Model	Task SLC			Task FLC		
	P	R	F ₁	P	R	F ₁
All-Propaganda	23.92	100.0	38.61	-	-	-
BERT	63.20	53.16	57.74	21.48	21.39	21.39
Joint	62.84	55.46	58.91	20.11	19.74	19.92
Granu	62.80	55.24	58.76	23.85	20.14	21.80
Multi-Granularity						
ReLU	60.41	61.58	60.98	23.98	20.33	21.82
Sigmoid	62.27	59.56	60.71	24.42	21.05	22.58

Table 1: Sentence-level (left) and fragment-level experiments (right). *All-propaganda* is a baseline that always output the propaganda class.

4 Experiments and Evaluation

We used the PyTorch⁴ framework and the pretrained BERT model,⁵ which we fine-tuned for our tasks.⁶ To deal with class imbalance, we give weight to the binary cross-entropy according to the proportion of positive samples. For the α in the joint loss function, we use 0.9 for sentence classification, and 0.1 for word-level classification. In order to reduce the effect of random fluctuations for BERT, all the reported numbers are the average of three experimental runs with different random seeds. As it is standard, we tune our models on the dev partition and we report results on the test partition.

The left side of Table 1 shows the performance for the three baselines and for our multi-granularity network on the FLC task. For the latter, we vary the degree to which the gate function is applied: using ReLU is more aggressive compared to using the Sigmoid, as the ReLU outputs zero for a negative input. Table 1 (right) shows that using additional information from the sentence-level for the token-level classification (BERT-Granularity) yields small improvements. The multi-granularity models outperform all baselines thanks to their higher precision. This shows the effect of the model excluding sentences that it determined to be non-propagandistic from being considered for token-level classification.

The right side of Table 1 shows the results for the SLC task. We apply our multi-granularity network model to the sentence-level classification task to see its effect on low granularity when we train the model with a high granularity task. Interestingly, it yields huge performance improvements on the sentence-level classification result. Compared to the BERT baseline, it increases the recall by 8.42%, resulting in a 3.24% increase of the F₁ score. In this case, the result of token-level classification is used as additional information for the sentence-level task, and it helps to find more positive samples. This shows the opposite effect of our model compared to the FLC task.

5 Conclusions

We have argued for a new way to study propaganda in news media: by focusing on identifying the instances of use of specific propaganda techniques. Going at this fine-grained level can yield more reliable systems and it also makes it possible to explain to the user why an article was judged as propagandistic by an automatic system.

We experimented with a number of BERT-based models and devised a novel architecture which outperforms standard BERT-based baselines. Our fine-grained task can complement document-level judgments, both to come out with an aggregated decision and to explain why a document—or an entire news outlet—has been flagged as potentially propagandistic by an automatic system.

In future work, we plan to include more media sources, especially from non-English-speaking media and regions. We further want to extend the tool to support other propaganda techniques.

⁴<http://pytorch.org>

⁵<http://github.com/huggingface/pytorch-pretrained-BERT>

⁶Our source code together with the dataset are available in GitHub: <http://anonymous.for.review>.

6 Acknowledgements

This research is part of the Propaganda Analysis Project,⁷ which is framed within the Tanbih project.⁸ The Tanbih project aims to limit the effect of “fake news”, propaganda, and media bias by making users aware of what they are reading, thus promoting media literacy and critical thinking. The project is developed in collaboration between the Qatar Computing Research Institute (QCRI), HBKU and the MIT Computer Science and Artificial Intelligence Laboratory (CSAIL).

References

- [1] A. Barrón-Cedeño, G. Da San Martino, I. Jaradat, and P. Nakov. Proppy: A system to unmask propaganda in online news. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence, AAAI '19*, pages 9847–9848, Honolulu, HI, USA, 2019.
- [2] A. Barrón-Cedeño, I. Jaradat, G. Da San Martino, and P. Nakov. Proppy: Organizing the news based on their propagandistic content. *Information Processing & Management*, 56(5):1849–1864, 2019.
- [3] G. Da San Martino, S. Yu, A. Barrón-Cedeño, R. Petrov, and P. Nakov. Fine-grained analysis of propaganda in news articles. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019*, pages 5640–5650, Hong Kong, China, 2019.
- [4] L. Dan. Techniques for the Translation of Advertising Slogans. In *Proceedings of the International Conference Literature, Discourse and Multicultural Dialogue, LDMD '15*, pages 13–23, Mures, Romania, 2015.
- [5] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. 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 '19*, pages 4171–4186, Minneapolis, MN, USA, 2019.
- [6] J. Goodwin. Accounting for the force of the appeal to authority. In *Proceedings of the 9th International Conference of the Ontario Society for the Study of Argumentation, OSSA '11*, pages 1–9, Ontario, Canada, 2011.
- [7] I. Habernal, H. Wachsmuth, I. Gurevych, and B. Stein. Before name-calling: Dynamics and triggers of ad hominem fallacies in web argumentation. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT '18*, pages 386–396, New Orleans, LA, USA, 2018.
- [8] B. D. Horne, S. Khedr, and S. Adali. Sampling the news producers: A large news and feature data set for the study of the complex media landscape. In *Proceedings of the Twelfth International AAAI Conference on Web and Social Media, ICWSM '18*, Stanford, CA, USA, 2018.
- [9] C. R. Miller. The Techniques of Propaganda. From “How to Detect and Analyze Propaganda,” an address given at Town Hall, 1939. The Center for learning.
- [10] M. Potthast, B. Stein, A. Barrón-Cedeño, and P. Rosso. An evaluation framework for plagiarism detection. In *Proceedings of the 23rd international conference on computational linguistics: Posters, COLING '10*, pages 997–1005, Beijing, China, 2010.
- [11] H. Rashkin, E. Choi, J. Y. Jang, S. Volkova, and Y. Choi. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP '17*, pages 2931–2937, Copenhagen, Denmark, 2017.
- [12] D. Walton. *The straw man fallacy*. Royal Netherlands Academy of Arts and Sciences, 1996.
- [13] A. Weston. *A rulebook for arguments*. Hackett Publishing, 2018.

⁷<http://propaganda.qcri.org>

⁸<http://tanbih.qcri.org>