

Claim Extraction From Unstructured Text

Majid Zarharan , Zahra Hosseini, Babak Behkamkia, Sauleh Eetemadi

April 2022

The spread of false information can lead to severe different problems [1]. Given that false information can be produced at a large scale and disseminated rapidly, it would be extremely difficult to detect and track them manually. Moreover, false claims are often wrapped in some true information. For example, we can consider a text that consists of one or more paragraphs, and all of the text is correct except a piece of information that may be a section of a sentence.

There are several works that analyze a sentence to determine whether it is a claim or not [2, 3, 4], but in our approach, even a sentence can encompass several independent claims. To cite some examples, we collected two claims from the FEVER [5] dataset with the NEI label.

The first example: Cameroon lists three definitely endangered languages, 13 severely endangered, and 16 critically endangered from among its at least 250 languages. There are seven different claims in this sentence:

1. Cameroon lists three definitely endangered languages.
2. Cameroon lists 13 severely endangered languages.
3. Cameroon lists 16 critically endangered languages.
4. Cameroon lists three definitely endangered languages from among its at least 250 languages.
5. Cameroon lists 13 severely endangered languages from among its at least 250 languages.
6. Cameroon lists 16 critically endangered languages from among its at least 250 languages.
7. Cameroon lists 250 languages.

The second example: Ryan Mathews of the 2003 Kansas State Wildcats football team posted a higher rushing average than his teammate Ell Roberson. There are three different claims in this sentence:

1. Ryan Mathews was a member of the Kansas State Wildcats football team.
2. Ryan Mathews posted a higher rushing average than his teammate Ell Roberson in 2003.
3. Ryan Mathews and Ell Roberson were teammates in 2003.

Using claim extraction, we can assess a text with any number of claims. Claim extraction is vital for the assessment of all existing claims in real-time news articles, public health texts, biomedical texts, etc. To begin, the text is divided into sentences. All possible claims are then extracted from each sentence and assessed individually. Lastly, we can conclude the label of the whole text by considering all extracted claims and their veracity. In addition, we show the veracity of the text in detail and plot the veracity of each piece of the text. For example:

Input text: Cristiano Ronaldo, who has scored over 800 senior career goals for club and country and has played for Real Madrid, has won two Ballon d’Or awards. In 2011, he lost his best friend.

At first, we break down the text into sentences, and there are two sentences in this example. For the next step, we extract all claims from these sentences which results in the following claims and assess¹ each claim individually:

1. Cristiano Ronaldo has scored over 800 senior career goals for club and country. (Support)
2. Cristiano Ronaldo has played for Real Madrid. (Support)
3. Cristiano Ronaldo has won two Ballon d’Or awards. (Refute)
4. In 2011, Cristiano Ronaldo lost his best friend. (Not Enough Info)

The final label for the text can be *Somewhat True* and we can plot the following result:

Cristiano Ronaldo, who has scored over 800 senior career goals for club and country and has played for Real Madrid, has won two Ballon d’Or awards. In 2011, he lost his best friend.

Note: In this work, we will want to focus on claim extraction.

References

- [1] W. Y. Wang, ““liar, liar pants on fire”: A new benchmark dataset for fake news detection,” in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*. Vancouver, Canada: Association for Computational Linguistics, Jul. 2017, pp. 422–426. [Online]. Available: <https://www.aclweb.org/anthology/P17-2067>
- [2] A. A. Prabhakar, S. Mohtaj, and S. Möller, “Claim extraction from text using transfer learning,” in *Proceedings of the 17th International Conference on Natural Language Processing (ICON)*. Indian Institute of Technology Patna, Patna, India: NLP Association of India (NLP AI), Dec. 2020, pp. 297–302. [Online]. Available: <https://aclanthology.org/2020.icon-main.39>
- [3] B. Adler and G. Boscaini-Gilroy, “Real-time claim detection from news articles and retrieval of semantically-similar factchecks,” *ArXiv*, vol. abs/1907.02030, 2019.
- [4] T. Jansen and T. Kuhn, “Extracting core claims from scientific articles,” in *BNAIC 2016: Artificial Intelligence*, ser. Communications in Computer and Information Science book series (CCIS), T. Bosse, Ed. Springer, 2017, pp. 32–46, in Post-proceedings of the 28th Benelux Conference on Artificial Intelligence (BNAIC 2016).
- [5] J. Thorne, A. Vlachos, C. Christodoulopoulos, and A. Mittal, “FEVER: a large-scale dataset for fact extraction and VERification,” in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. New Orleans, Louisiana: Association for Computational Linguistics, Jun. 2018, pp. 809–819. [Online]. Available: <https://www.aclweb.org/anthology/N18-1074>

¹we assessed these claims by considering this [link](#).