Abstract

This poster presents a method to identify the hidden structures of documents. Each document has a rhetorical structure, which expresses the relations among clauses. Since building a rhetorical structure parser is based on supervised learning, it requires large amounts of manually annotated training data for accurate parsing. However, conventional methods suffer from a lack of training data, resulting in poor performance because manual annotation is guite labor intensive. To tackle this problem, we propose a method that uses silver data: automatically annotated pseudo-labeled data. We pre-trained the parser with silver data and fine-tuned it with gold data: manually annotated data. Our experimental results demonstrated that our method achieved the best performance. The new parser will contribute to various natural language processing applications, such as machine translation and automatic summarization.



[1] N. Koabayashi, T. Hirao, H. Kamigaito, M. Okumura, M. Nagata, "Improving Neural RST Parsing Model with Silver Agreement Subtrees," in Proc. 2021 Annual Conference of the Noth American Chapter of the Association for Computational Linguistics, 2021.

Contact

Tsutomu Hirao / Linguistic Intelligence Research Group, Innovative Communication Laboratory Email: cs-openhouse-ml@hco.ntt.co.jp