

Minghui Meng, Yuzhuo Wang, Chengzhi Zhang

Department of Information Management, Nanjing University of Science and Technology, Nanjing, China

Introduction

Peer reviews for academic articles are domain experts' comments on the paper submitted to a journal or conference, which reflect the overall impression of the reviewer as well as the detailed comments. By mining peer reviews at aspect level, we can find the aspects concerned by reviewers, which can provide direction for inexperienced paper submitters to optimize their writing. Some scholars have studied it at aspect level, but the aspects are coarse granularity and lack of more detailed comments. Moreover, reviewers of different disciplines may focus on different aspects. Therefore, we want to find more fine-grained multi-level aspects from multi-disciplinary peer reviews, including domain independent aspects and domain related aspects. We take peer reviews in Nature Communications as an example and propose a general method for building multi-level aspects of peer reviews for academic articles.

Method

Data collection

Nature Communications (NC) is a multidisciplinary journal, which has 5 first-level disciplines and 71 second-level disciplines. We regard "NC" as the zero-level discipline. We collect 187,971 peer review documents of articles published from 2016 to 2020, details of the NC's first-level disciplines are shown in table 1:

Table 1. Distribution of NC Peer Reviews

Discipline	Papers	Reviews
Biological sciences	21,683	98,585
Earth and environmental sciences	2,588	11,632
Health sciences	6,090	28,674
Physical sciences	10,180	45,222
Scientific community and society	832	3,858

Candidate aspects extraction

We use the Double Propagation algorithm (Qiu et al., 2011) to extract aspects from peer reviews. First, we employ the StanfordNLP tool to tag the part of speech of peer reviews and parse sentences. Second, we select words in opinion lexicon from Liu as seed opinions. Next, we consider aspects to be nouns and opinions to be adjectives and limit the dependency relations between aspects and opinions to mod, subj, etc. Finally, we select all the words matching the above rules as the candidate aspects.

Multi-level aspects determination

We divide multi-level aspects into domain-independent common aspects (DICA), domain-related common aspects (DRCA), and domain-related special aspects (DRSA). DICA is the aspect that often appears and distributes evenly in every first-level domains. DRCA often appears in a certain first-level disciplinary domain and is evenly distributed in its sub-domains. DRSA is the aspect unevenly distributed in each second-level domains and often appears in a certain second-level domain. Multi-level aspects are determined based on the number of reviews, the distribution uniformity, and particularity of candidate aspects at all levels of domain peer reviews. We use the inter-domain entropy (IDE) and termhood (Chang, 2005) to measure the distribution uniformity and particularity of aspects in various domains, the calculation formulas are as follows:

$$(1) IDE(w_i) = -\sum_j P_{ij} \log P_{ij} \quad (2) Termhood(w_{ij}) = \frac{1}{n_{ij}} \times \log_2 \left[\frac{N}{Nd_i} \right]$$

Where P_{ij} is the probability of candidate words w_i in domain J , $Nd_i = 2^{IDE(w_i)}$, N is the total number of the disciplinary domains.

Aspects clustering

We select the Affinity Propagation (AP) algorithm to cluster aspects. We train word2vec model to represent word vector, then use cosine similarity to calculate the similarity between two aspects and choose the Silhouette Coefficient (SC) to evaluate the clustering result.

Results

The extraction result of candidate aspects

We respectively extract the candidate aspects from NC's every-level discipline peer reviews. As a result, we extract 5896 zero-level candidate aspects, 5 groups of first-level candidate aspects, and 71 groups of second-level candidate aspects.

Table 2. Multi-Level Aspects

Level	Category	Domain	Aspects
0	DICA	NC	Data\result\method\figure\analysis\change\experiment\effect\text\level\model\topic\example\...\discussion\
1	DICA	BIO	Assay\target\domain\...\specificity
2	DRSA	CBB	Object\update\classifier\...\paradigm

The determination result of multi-level aspects

In NC peer reviews, we obtain 598 DICA, 5 groups of DRCA, and 71 groups of DRSA. Table2 is the partial result of the NC multi-level aspects, including NC's DICA, Biological sciences (BIO)'s DRCA and Computational biology and bioinformatics (CBB)'s DRSA.

The clustering result of zero-level discipline aspects

We cluster NC's DICA, when preference = -40, damping = 0.65, the SC value was the largest, we regard it as the optimal clustering result. Table3 is part of the clustering results of NC's DICA.

Table 3. The NC's DICA Clusters

Cluster	Aspects	#Member
Result	Result\conclusion\finding\...\claim	19
Impact	Impact\lack\contribution\...\focus	48
Method	Method\approach\technique\...\strategy	19
Figure	Figure\image\table\legend\...\panel	29

Conclusion

This paper proposes a method for building multi-level aspects of peer reviews for academic articles, which aims to comprehensively extract the multi-level aspects in peer reviews. The multi-level aspects reflect the key aspects that reviewers focus on, such as result and method, which can provide reference for inexperienced submitters to optimize their research design and writing.

References

- Chakraborty, S., Goyal, P., Mukherjee, A. (2020). Aspect-based Sentiment Analysis of Scientific Reviews. Proceedings of the ACM/IEEE Joint Conference on Digital Libraries in 2020(JCDL '20) (pp.63-81). New York, USA: ACM.
- Qiu, G., Liu, B., Bu, J., Chen, C. (2011). Opinion word expansion and target extraction through double propagation. Computational Linguistics, 37(1), 9-27.
- Chang J. (2005). Domain Specific Word Extraction from Hierarchical Web Documents: a First Step toward Building Lexicon Trees from Web Corpora. Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing (pp.64-71). Stroudsburg, PA: ACL.
- Frey, B J, & Dueck, D. (2007). Clustering by passing messages between data points. Science, 315(5814), 972-976.

Contact

Minghui Meng (mengmh@njust.edu.cn)

Chengzhi Zhang (Corresponding author, zhangcz@njust.edu.cn)