| cases | doc_1 | | doc_2 | | decision | id |
|---|---|---|---|---|---|---|
| | authors | • Xintao Chu<br>• Jianping Liu<br>• Jian Wang<br>• Xiaofeng Wang<br>• Yingfei Wang<br>• Meng Wang<br>• Xunxun Gu | authors | • Xintao Chu<br>• Jianping Liu<br>• Jian Wang<br>• Xiaofeng Wang<br>• Yingfei Wang<br>• Meng Wang<br>• Xunxun Gu | DUPLICATES | 6 |
| | title | CSDR-BERT: a pre-trained scientific dataset match model for Chinese Scientific Dataset Retrieval | title | CSDR-BERT: a pre-trained scientific dataset match model for Chinese Scientific Dataset Retrieval | | |
| | publication_date | 2023-01-30 07:12:38+00:00 | publication_date | 2023-01-31 00:00:00 | | |
| | source | SupportedSources.ARXIV | source | SupportedSources.INTERNET_ARCHIVE | | |
| | journal | None | journal | | | |
| | volume | | volume | | | |
| | doi | | doi | | | |
| | urls | • http://arxiv.org/pdf/2301.12700v3<br>• http://arxiv.org/abs/2301.12700v3<br>• http://arxiv.org/pdf/2301.12700v3 | urls | • https://web.archive.org/web/20230206013603/https://arxiv.org/pdf/2301.12700v2.pdf | | |
| | id | id-565247275900382637 | id | id4508891993084394901 | | |
| | abstract | As the number of open and shared scientific datasets on the Internet increases under the open science movement, efficiently retrieving these datasets is a crucial task in information retrieval (IR) research. In recent years, the development of large models, particularly the pre-training and fine-tuning paradigm, which involves pre-training on large models and fine-tuning on downstream tasks, has provided new solutions for IR match tasks. In this study, we use the original BERT token in the embedding layer, improve the Sentence-BERT model structure in the model layer by introducing the SimCSE and K-Nearest Neighbors method, and use the cosent loss function in the optimization phase to optimize the target output. Our experimental results show that our model outperforms other competing models on both public and self-built datasets through comparative experiments and ablation implementations. This study explores and validates the feasibility and efficiency of pre-training techniques for semantic retrieval of Chinese scientific datasets. | abstract | As the number of open and shared scientific datasets on the Internet increases under the open science movement, efficiently retrieving these datasets is a crucial task in information retrieval (IR) research. In recent years, the development of large models, particularly the pre-training and fine-tuning paradigm, which involves pre-training on large models and fine-tuning on downstream tasks, has provided new solutions for IR match tasks. In this study, we use the original BERT token in the embedding layer, improve the Sentence-BERT model structure in the model layer by introducing the SimCSE and K-Nearest Neighbors method, and use the cosent loss function in the optimization phase to optimize the target output. Our experimental results show that our model outperforms other competing models on both public and self-built datasets through comparative experiments and ablation implementations. This study explores and validates the feasibility and efficiency of pre-training techniques for semantic retrieval of Chinese scientific datasets. | | |
| | versions | | versions | | | |