| cases | | doc_1 | | doc_2 | decision | id |
|---|---|---|---|---|---|---|
| | authors | <ul><li>Tianyu Gao</li><li>Xingcheng Yao</li><li>Danqi Chen</li></ul> | authors | <ul><li>Tianyu Gao</li><li>Xing-Cheng Yao</li><li>Danqi Chen</li></ul> | DUPLICATES | 156 |
| | title | SimCSE: Simple Contrastive Learning of Sentence Embeddings | title | SimCSE: Simple Contrastive Learning of Sentence Embeddings | | |
| | publication_date | 2022-05-18 00:00:00 | publication_date | 2021-04-18 00:00:00 | | |
| | source | SupportedSources.INTERNET_ARCHIVE | source | SupportedSources.OPENALEX | | |
| | journal | | journal | arXiv (Cornell University) | | |
| | volume | | volume | | | |
| | doi | | doi | 10.48550/arxiv.2104.08821 | | |
| | urls | <ul><li>https://web.archive.org/web/20220526183956/https://arxiv.org/pdf/2104.08821v4.pdf</li></ul> | urls | <ul><li>https://openalex.org/W3213189520</li><li>https://doi.org/10.48550/arxiv.2104.08821</li><li>http://arxiv.org/pdf/2104.08821</li></ul> | | |
| | id | id9115064052159856573 | id | id-7788975266327500169 | | |
| | abstract | This paper presents SimCSE, a simple contrastive learning framework that greatly advances state-of-the-art sentence embeddings. We first describe an unsupervised approach, which takes an input sentence and predicts itself in a contrastive objective, with only standard dropout used as noise. This simple method works surprisingly well, performing on par with previous supervised counterparts. We find that dropout acts as minimal data augmentation, and removing it leads to a representation collapse. Then, we propose a supervised approach, which incorporates annotated pairs from natural language inference datasets into our contrastive learning framework by using "entailment" pairs as positives and "contradiction" pairs as hard negatives. We evaluate SimCSE on standard semantic textual similarity (STS) tasks, and our unsupervised and supervised models using BERT base achieve an average of 76.3% and 81.6% Spearman's correlation respectively, a 4.2% and 2.2% improvement compared to the previous best results. We also show -- both theoretically and empirically -- that the contrastive learning objective regularizes pre-trained embeddings' anisotropic space to be more uniform, and it better aligns positive pairs when supervised signals are available. | abstract | | | |
| | versions | | versions | | | |