| | doc_1 | doc_2 | decision | id |
|---|---|---|---|---|
| cases | **authors**: • Lee, Jaejun • Tang, Raphael • Lin, Jimmy<br><br>**title**: What Would Elsa Do? Freezing Layers During Transformer Fine-Tuning<br>**publication_date**: 2019-11-08 00:00:00<br>**source**: SupportedSources.OPENALEX<br>**journal**: arXiv (Cornell University)<br>**volume**:<br>**doi**: 10.48550/arxiv.1911.03090<br>**urls**: • https://openalex.org/W4288026527 • https://doi.org/10.48550/arxiv.1911.03090 • http://arxiv.org/pdf/1911.03090<br>**id**: id5285092138517909115<br>**abstract**:<br>**versions**: | **authors**: • Jaejun Lee • Raphael Tang • Jimmy Lin<br><br>**title**: What Would Elsa Do? Freezing Layers During Transformer Fine-Tuning<br>**publication_date**: 2019-11-08 07:05:20+00:00<br>**source**: SupportedSources.ARXIV<br>**journal**: None<br>**volume**:<br>**doi**:<br>**urls**: • http://arxiv.org/pdf/1911.03090v1 • http://arxiv.org/abs/1911.03090v1 • http://arxiv.org/pdf/1911.03090v1<br>**id**: id-2980461387407644797<br>**abstract**: Pretrained transformer-based language models have achieved state of the art across countless tasks in natural language processing. These models are highly expressive, comprising at least a hundred million parameters and a dozen layers. Recent evidence suggests that only a few of the final layers need to be fine-tuned for high quality on downstream tasks. Naturally, a subsequent research question is, "how many of the last layers do we need to fine-tune?" In this paper, we precisely answer this question. We examine two recent pretrained language models, BERT and RoBERTa, across standard tasks in textual entailment, semantic similarity, sentiment analysis, and linguistic acceptability. We vary the number of final layers that are fine-tuned, then study the resulting change in task-specific effectiveness. We show that only a fourth of the final layers need to be fine-tuned to achieve 90% of the original quality. Surprisingly, we also find that fine-tuning all layers does not always help.<br>**versions**: | DUPLICATES | 339 |