

	doc_1		doc_2		decision	id
cases			authors	<ul style="list-style-type: none"> Jinwoo Shin Jaeho Lee Kimin Lee Sangwoo Mo Seung Jun Moon 	DUPLICATES	268
	authors	<ul style="list-style-type: none"> Seung Ki Moon Sangwoo Mo Kimin Lee Jae-Ho Lee Jinwoo Shin 	title	MASKER: Masked Keyword Regularization for Reliable Text Classification		
			publication_date	2020-12-17 00:00:00		
			source	SupportedSources.PAPERS_WITH_CODE		
			journal			
			volume			
			doi			
			urls	<ul style="list-style-type: none"> https://arxiv.org/pdf/2012.09392v1.pdf https://github.com/alinlab/MASKER 		
			id	id1670891010620103543		
			abstract	Pre-trained language models have achieved state-of-the-art accuracies on various text classification tasks, e.g., sentiment analysis, natural language inference, and semantic textual similarity. However, the reliability of the fine-tuned text classifiers is an often overlooked performance criterion. For instance, one may desire a model that can detect out-of-distribution (OOD) samples (drawn far from training distribution) or be robust against domain shifts. We claim that one central obstacle to the reliability is the over-reliance of the model on a limited number of keywords, instead of looking at the whole context. In particular, we find that (a) OOD samples often contain in-distribution keywords, while (b) cross-domain samples may not always contain keywords; over-relying on the keywords can be problematic for both cases. In light of this observation, we propose a simple yet effective fine-tuning method, coined masked keyword regularization (MASKER), that facilitates context-based prediction. MASKER regularizes the model to reconstruct the keywords from the rest of the words and make low-confidence predictions without enough context. When applied to various pre-trained language models (e.g., BERT, RoBERTa, and ALBERT), we demonstrate that MASKER improves OOD detection and cross-domain generalization without degrading classification accuracy. Code is available at https://github.com/alinlab/MASKER.		
			versions			