

cases	doc_1		doc_2				decision	id
							DUPLICATES	117
	authors	<ul style="list-style-type: none">Alexander MeinkeJulian BitterwolfMatthias Hein	authors	<ul style="list-style-type: none">Matthias HeinJulian BitterwolfAlexander Meinke				
	title	Provably Robust Detection of Out-of-distribution Data (almost) for free.	title	Provably Robust Detection of Out-of-distribution Data (almost) for free				
	publication_date	2021-06-08 00:00:00	publication_date	2021-06-08 00:00:00				
	source	SupportedSources.OPENALEX	source	SupportedSources.PAPERS_WITH_CODE				
	journal	arXiv (Cornell University)	journal					
	volume		volume					
	doi	None	doi					
	urls	<ul style="list-style-type: none">https://openalex.org/W3172596993	urls	<ul style="list-style-type: none">https://arxiv.org/pdf/2106.04260v2.pdfhttps://github.com/AlexMeinke/Provable-OOD-Detectionhttps://openreview.net/pdf?id=qDx6DXD3Fzt				
	id	id-264729082788196000	id	id-9201529456528127826				
	abstract		abstract	The application of machine learning in safety-critical systems requires a reliable assessment of uncertainty. However, deep neural networks are known to produce highly overconfident predictions on out-of-distribution (OOD) data. Even if trained to be non-confident on OOD data, one can still adversarially manipulate OOD data so that the classifier again assigns high confidence to the manipulated samples. We show that two previously published defenses can be broken by better adapted attacks, highlighting the importance of robustness guarantees around OOD data. Since the existing method for this task is hard to train and significantly limits accuracy, we construct a classifier that can simultaneously achieve provably adversarially robust OOD detection and high clean accuracy. Moreover, by slightly modifying the classifier's architecture our method provably avoids the asymptotic overconfidence problem of standard neural networks. We provide code for all our experiments.				
	versions		versions					