

URL: <https://stvp.stanford.edu/clips/reducing-bias-in-datasets>

At least for now, observes Scale AI founder and CEO Alexandr Wang, we don't have an easy solution that can completely eliminate data bias. He emphasizes the essential role high-quality data plays in building responsible, unbiased AI systems, and explores how Scale works with customers to proactively reduce data bias.



Transcript

Alexandr To have responsibility, 00:00:03,810 it's incredibly important to build high quality and representative data.. You know, I actually think it's critical, because at its core, this thing that I talked about earlier is that, you know, the data is almost like the food that you feed to these machine learning systems and ultimately you are what you eat.. And so, the sort of the, data is sort of the, one of the fundamental areas where these biases or these kinds of really critical real-world issues can sort of stem from.. And, really the way we think of it, I was like, how do we, how do we build technology? How do we build tools? How do we build things to enable either our customers or ourselves to ensure minimal bias and minimal harmful results with these machines, with the datasets that's part of the machine learning systems, and or, how do we build systems that are going to flag that, or identify when those biases may exist as proactively as possible.. So that then we know that we have to go fix that problem long before it ends up in the hands of a consumer or some sort of critical decision making progress.. So for example, we recently worked with like a bunch of medical researchers on automated medical imaging analysis.. You know, the sort of exact problem we worked with the MIT media lab and analyzed this exact problem that we mentioned of tons and tons of clinical images and found there were a lot more light-skinned images and dark-skinned images.. And what we did is we, actually basically helped them replenish or sort of debias the dataset, and add more data to the unrepresentative classes, either with real-world data or using data augmentation or synthetic data.. And we're able to significantly debias the output or the outcomes from the machine learning algorithm.. And so I think, you know, it's not an easy problem in any way, you know, this is actually a sort of like, you know, this is one of the (mumbles) hard problems of AI, in the sense that there's no easy solution..

It's not like, you know, I don't expect that next year we're going to get some, you know, crazy machine learning architecture that all of a sudden solves the problem of data bias or machine learning bias.. It really is all about the nitty-gritty details about what is your dataset look like.. How do you know ahead of time when there's bias in the dataset, and then how do you fix that? Or how do you practically resolve that? And how do you keep going through that process?..