

This article contains excerpts from the New York Times "Dealing with bias in artificial intelligence" from November 2019. Specifically the condensed comments of Dr. Koller, an adjunct professor in the computer science department at Stanford University.

1 You could mean bias in the sense of racial bias, gender bias. For example, you do a search for C.E.O. on Google Images, and up come 50 images
2 of white males and one image of C.E.O. Barbie. That's one aspect of bias.

3 Another notion of bias, one that is highly relevant to my work, are cases in which an algorithm is latching onto something that is meaningless
4 and could potentially give you very poor results. For example. imagine that you are trying to predict fractures from X-ray images in data from
5 multiple hospitals. If you are not careful, the algorithm will learn to recognize which hospital generated the image. Some X-ray machines have
6 different characteristics in the image they produce than other machines, and some hospitals have a much larger percentage of fractures than
7 others. And so, you could actually learn to predict fractures pretty well on the data set that you were given simply by recognizing which
8 hospital did the scan, without actually looking at the bone. The algorithm is doing something that appears to be good but it is actually doing
9 it for the wrong reasons. The causes are the same in the sense that these are all about how the algorithm latches onto things that it shouldn't
10 latch onto in making its prediction.

11 To recognize and address these situations, you have to make sure that you test the algorithm in a regime
12 that is similar to how it will be used in the real world. So, if your machine-learning algorithm is one that is
13 trained on the data from a given set of hospitals, and you will only use it in those same set of hospitals, then
14 latching onto which hospital did the scan could well be a reasonable approach. It's effectively letting the
15 algorithm incorporate prior knowledge about the patient population in different hospitals. The problem
16 really arises if you're going to use that algorithm in the context of another hospital that wasn't in your data
17 set to begin with. Then, you're asking the algorithm to use these biases that it learned on the hospitals that
18 it trained on, on a hospital where the biases might be completely wrong.

19 Over all, there's not nearly as much sophistication as there needs to be out there for the level of rigor that
20 we need in terms of the application of data science to real-world data, and especially biomedical data.

