

# Data Quality: The Risks Of Dirty Data And AI



Intel AI BRANDVOICE | Paid Program



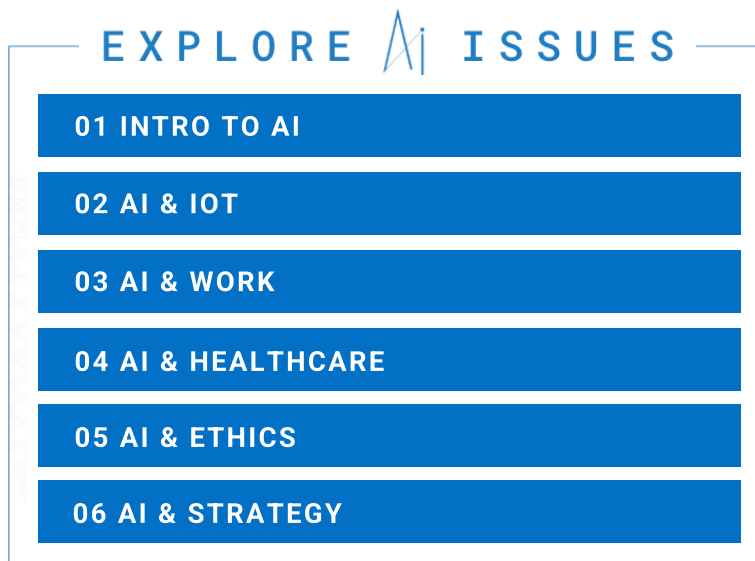
Mar 27, 2019, 01:21pm EDT

Updated Jun 5, 2019, 04:36pm EDT

*By Jason Compton*

**A**n artificial intelligence (AI) is only as useful as the data used to train it. Training AI with dubious data leads to questionable decisions.

An AI project should thus give a thorough airing to ethical considerations from the very second it launches—since the assumptions, experiences, prejudices and mental habits of the human beings who create each input, and interpret each output, will inevitably inform it.



“AI is a mirror that reflects back to us the bias that already exists in our society,” [writes](#) Kathy Baxter, architect of ethical AI practice at Salesforce.

## Navigating The Difficulties Of Dirty Data

Skewed data sets are a problem as old as human cognition. It is well understood that people with strong positive or negative opinions are more likely to offer commentary or respond to a survey than those in the muddled middle. And data distortion is a widespread problem—whether or not an organization uses AI.

According to a 2017 Harvard Business Review [study](#), almost every company has data quality shortfalls. Nearly half have data quality issues entailing clear and negative business consequences.

Data biases and data hygiene problems can have short- and long-term consequences—

consequences that are not necessarily evenly distributed.

Anna Bethke, the head of AI for social good at Intel's Artificial Intelligence Products Group, provides some simple examples of how data quality and data hygiene can make a difference. AI trained on a relatively unbiased data set would be equally able to classify a Western, Eastern or African bride, whereas one trained almost exclusively on images from Western cultures might misidentify brides wearing anything other than white gowns. And the datasets on which we often train language models are able to correctly ascertain that “man is to king as woman is to queen”—but can also infer that “[man is to computer programmer as woman is to homemaker](#).”

The problem isn't the algorithm, but the dirty data fed into it. That data is, in turn, the product of conscious or unconscious biases in our society.

In the short term, such dirty data problems might visibly harm only the group directly connected to the AI's mistaken conclusion. Over the long term, however, there could be reputational damage, and the business in question could miss out on opportunities due to an accretion of poor decision-making.

*[Learn more about how companies are leveraging AI today.](#)*

## **Grappling With Ethics And Fairness**

The AI Now Institute recently published an [assertive list](#) of some of the lowlights in AI ethics over the past year, along with a [report](#) highlighting the growing ethical risks in surveillance. The latter manifest themselves both in the form of unchecked capacity to monitor individuals beyond their ability to reasonably consent and in the controversial practice of assessing affect (mood, body language and tendencies) through algorithms.

There is no single, simple answer to the ethics problem. A recent survey on the topic [found](#) that more than 20 competing concepts of fairness relating to AI and machine learning had been published in the past decade.

“There’s no clear, definitive way to just solve for one variable, or to just use a particular algorithm, the way we solve other technology problems,” Bethke says. “This is also a social science problem.”

And it’s a big one—it’s been a topic of serious consideration by the largest technology industry bodies for decades.

The Association for Computing Machinery’s Code of Ethics [states the case](#) plainly: “Extraordinary care should be taken to identify and mitigate potential risks in machine learning systems. A system for which future risks cannot be reliably predicted requires frequent reassessment of risk as the system evolves in use, or it should not be deployed.”

The risks are greater than those that attend simply taking a new system live or performing a server migration, and AI professionals should treat them accordingly.

## Tackling The Basics Of Data Quality

In fact, ethical issues at the intersection of data quality, data hygiene and technology have been [openly discussed](#) since the 1960s, and they remain a hot topic for debate. Complex as this stuff is, it's actually relatively easy to get started on the fundamentals.

Here are several ways to take on (and head off) the ethical challenges that can creep into any AI project.

### 1. Recognize that “bias” doesn’t mean “malice.”

Start with the reality that bias itself can be a loaded term and may make people reflexively defensive, protesting that they are not personally biased or prejudiced against a particular group. But algorithms can reflect the unconscious assumptions and blind spots of their creators, thereby amplifying a lack of experience or a narrow worldview.

And this doesn't apply only to algorithms. Data sets can also be invisibly biased.

For example, some facial recognition algorithms were trained on repositories of celebrity photographs. On the surface, this is a sensible design choice: Celebrity photos are

widely available, with a variety of high-quality images captured by skilled photographers.

But an algorithm trained only on celebrity appearances will not reflect the full diversity of the world population, which is less white and less male than the set of celebrities—and which, moreover, is not professionally styled for high-profile photo shoots.

A charge of bias is not a charge of malice. AI researchers and business users alike need to be ready to push past the urge to qualify or defend their efforts, and instead look at the facts of the matter.

## 2. Talk about it.

Industry is stepping up to the challenge, with interdisciplinary ethics panels, workshops and expert groups at every reputable machine learning event and ethical watchdogs in every conscientious AI laboratory.

But it's important to expand the data quality and data hygiene conversation beyond AI practitioners, data scientists and business owners.

“I recently heard a powerful piece of advice: When you build a system, talk to your neighbors, talk to a priest, talk to a 5-year old,” Bethke says. “Ask what they think about it.”

## 3. Question your results.

One of the clearest ways to trace the presence of bias is to audit the algorithm and look at why it's reaching its decisions. In addition to exposing dirty data problems or another data quality issue, you may learn that the algorithm has learned to look for the wrong things altogether.

As Bethke puts it, offering an example of this type of audit: “Did it classify an image as a dog because there’s actually a dog in the picture? Or was it focusing on the surroundings?”



#### 4. Use the available tools.

There's a growing range of practices and procedures with which to audit data quality and AI and data fairness. The free and adaptable [Ethics & Algorithms Toolkit](#) provides a process that professionals can use to evaluate the strengths, risks and biases of their algorithms. And [deon](#) provides an automated way to add an ethical checklist to data science projects.

5. Realize it's not either/or.

AI can still [accomplish good](#) even as discussions about fairness, inclusion and ethics continue.

“It’s important to be nondefensive about your AI project—take these ethical ideas, synthesize them and come up with a solution that’s beneficial for as many as possible,” Bethke explains.

*Jason Compton is a writer and reporter with extensive experience in enterprise tech. He is the former executive editor of CRM Magazine.*

[Learn more about how companies are leveraging AI today.](#)

Credit: Hero Images/GettyImages



**Intel AI**

You may know us for our processors. But we do so much more.  
Intel invents at the... **Read More**

Editorial Standards

Forbes Accolades

---