# Teaching Data Ethics on the Bleeding Edge

Tian Zheng, Isabelle Zaugg, Jonathan Reeve

From the spread of disinformation via social media, to class-biased dynamic pricing, to racial profiling in online systems that lead to "real-world" harms, teaching data ethics has never been more urgently needed in the data science curriculum. This article conducts a thorough review of the literature on data ethics and explores current trends among syllabi on the topic. It analyzes the overlaps and divides between various approaches to data ethics, from FAT (Fairness, Accountability, Transparency), to the Public Interest Technology movement, to calls to reimagine "digital justice" from Critical Race and Digital Studies scholars such as Ruha Benjamin, to the way classical philosophical texts on ethics from different cultural contexts are being applied to the digital age. Besides a traditional literature review, we also present a semantic, linked open data graph describing the relations between texts, courses, professors, and universities involved in teaching data ethics. Finally, we propose next steps in defining a cutting-edge approach to teaching data ethics that takes into account diverse understandings of the topic within the urgency of the moment. This approach will not only familiarize students with data ethics theories, but will push them to recognize a horizon of possible solutions and build a foundation of computational literacy to explore solutions in practice.

The semantic web, also known as "Web 3.0" or linked open data, is a relatively new system of conventions for standardizing and encoding graph data, such that it is universally interoperable, in a language known as RDF, or the Resource Description Framework. Some of the most well-known projects in the field include DBPedia, the set of parsed and inferred data from Wikipedia, and Wikidata, the data set which proposes to be the knowledge basis for Wikipedia. At its most basic, RDF data may be represented as a series of subject-verb-object triples, where each node has a stable URL. Social relationships between people, for instance, may be described as `<Bob> <is friends with> <Alice>`, where the angle-bracketed entities resolve to URIs. There exist a number of *ontologies*, or pre-defined sets of relations, which may be used to describe entities within their domains. For instance, in the social network example, the Friend-of-a-Friend (FOAF) ontology may be used to describe relationships between people. We use a number of ontologies in conjunction: the Curriculum Course Syllabus Ontology (CCSO) describes relations between courses, universities, syllabi, professors, and learning materials such as texts; the Bibliographic Ontology (Bibliontology) describes metadata for articles, books, videos, and other media; and the Citation

Typing Ontology (CiTO) describes citation relations between texts. Fig. 1 shows an example directed graph visualization, illustrating relations between these entities.

We collect data in a quasi-automated fashion, often beginning with course lists, such as the tech ethics curriculum list provided by Casey Fiesler et al. From there, given a course syllabus URL, we are able to automatically extract bibliographic references, and resolve them to stable identifiers at a number of bibliographic databases, such as CrossRef. These databases allow us to derive further information, such as the sources of funding for projects associated with publications. Universities and departments we then resolve to their Wikidata identifiers, which allows us to retrieve a considerable amount of additional information with which we can organize our data: for instance, the geographic coordinates of the university, and its date of foundation. We resolve instructors and authors to their ORCIDs, which allows us to retrieve other publications from the same author, and their past and present institutional affiliations. All of this all allows us to answer questions such as:

- What are the most-cited books and articles in the field of data ethics?
- Which are the books and articles most assigned in courses?
- What books are only cited in one geographic region (e.g., California), but nowhere else?
- What courses are cross-listed in the most number of departments?

We are currently building a website to visualize these connections, as a force-directed network visualization in JavaScript, so that it may be explored by a wider user base. This semantic web approach also allows us to be multilingual by default, since much of this data, such as that gleaned from Wikipedia, is available in many languages.

One next step for this project may include building a mechanism for users to submit their own data ethics courses to our database: this way, our literature review will always stay up to date. A further step will be to generalize this framework, so that it may be used to map any academic discipline: given a list of courses and their syllabus URLs, it will generate a map of all the academic entities involved: syllabi, required texts, instructors, and more.

As data ethicists ourselves, we care about openness and transparency, and so we have open-sourced this data, so that other researchers can use our work to answer their own questions. We hope that our framework may be used to help map the institutional knowledge structures of even more disciplines.
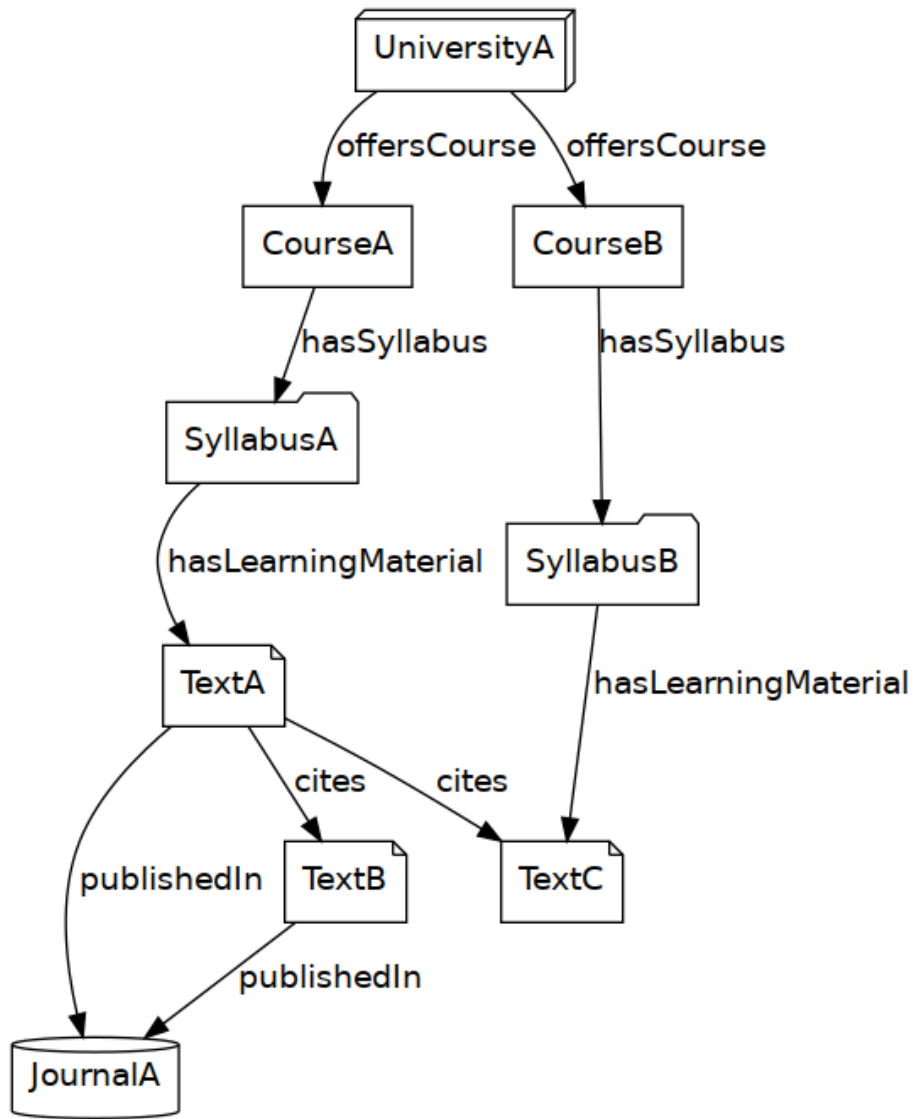
Figure 1: an example graph