

## Erasure, Misrepresentation and Confusion: Investigating JSTOR Topics on Women's and Race Histories

Sharon Block <sbblock\_at\_uci\_dot\_edu>, UC-Irvine

### Abstract

This article investigates the topic labeling system of a widely used full-text academic publication database, JSTOR, particularly in reference to colonial North American history scholarship. Using insights developed by critical algorithm and critical archival studies, it analyzes how JSTOR's topics repeatedly misrepresent and erase work in women's, African diasporic/African American, and Native American and settler colonial histories. The article discusses concerns over the power of metadata, the need for transparent and domain-expert-involved indexing processes, and digital providers' responsibilities to accurately categorize scholars' work. It particularly focuses on the potentially disproportionate harm done to traditionally marginalized fields of study through seemingly racist or sexist topical labeling that impedes knowledge discovery.

### Introduction

For at least two decades, scholars have written on the degree to which we live in an algorithmic culture, a computational theocracy, or been beholden to the power of computer algorithms [Granieri 2014] [Bogost 2015] [Introna and Nissenbaum 2000]. In recent years, scholars have published book-length critiques of the sexism and racism behind increasingly omnipresent and opaque search systems and mobile apps [Eubanks 2018] [Noble 2018] [O'Neil 2016] [Wachter-Boettcher 2017]. I build on these scholars' groundbreaking works on the hazards of algorithmic bias by analyzing one academic database's topical indexing functions. Beyond a critique of inaccuracies and omissions, I detail how such subject miscategorizations reinforce sexist and racist belief systems, thus having a disproportionate effect on marginalized groups and research.

1

In April 2018, I attempted to use the topic indexing system in the academic database, JSTOR, to prepare a state of the field presentation on early American women's history. I quickly realized that in several areas of my expertise (colonial North America, women's, race, Indigenous and African American histories), JSTOR's topic categorizations displayed concerning shortcomings. Articles that were focused entirely on women's history did not seem to be categorized by the topic of "women". Instead, some were mischaracterized with "men" as the most relevant topic. JSTOR's topics for African, African American, Native American and race histories showed misapplications and erasures as well, fundamentally distorting the content of scholarship in these fields.

2

Most scholars are at least passingly familiar with the Library of Congress Classification System and the controlled vocabulary (a set list of terms used for indexing and information retrieval) of the Library of Congress Subject Headings [Library of Congress Classification] [Library of Congress Subject and Genre/Form Headings]. These systems have structured knowledge for well over a century. Such categorization schemas have always been less than objective, as the 2016 struggle between the Library of Congress and House of Representatives over the subject heading "Illegal aliens" made abundantly clear [Peet 2016]. Scholars have pointed to anti-LGBT bias, Eurocentric and anti-Afrocentric biases, outdated terminology, and the limitations of discipline-based, hierarchical structure within the Library of Congress classification systems for many decades [Bethel 1994] [Christensen 2008] [Diao and Cao 2016] [Drabinski 2009] [Howard and Knowlton 2018]. Thus, concern over bias in indexing and classification schemas is not a recent

3

phenomenon.

However, the rise of digital databases and accompanying machine learning technologies has brought new concerns. Critical algorithm studies have developed in response to the influence of unknowable algorithms on consequential decisions and actions. At their most basic, algorithms are a list of programmed instructions, and can be millions of lines of proprietary coding, never seen or understood by end users. Scholars have called attention to the difficulty in “deconstructing the black boxes of Big Data” that create our algorithmic culture [Pasquale 2015, 6]. Proprietary algorithms that produce racist results, in particular, have been a repeated concern among data scientists and social justice advocates. [Buolamwini 2016] [Paul 2016] [Schwartzapfel 2019] [Ulloa 2019]. One of the most influential media and library information scholars, Safiya Noble, has convincingly argued that we must interrogate the “implications of the artificial intelligentsia for people who are already systematically marginalized and oppressed” [Noble 2018, 3]. Noble’s work documenting and challenging the racism reproduced by search engines (primarily Google) has been particularly impactful both within and outside academic discourse.

Digital Humanities scholars have argued that seemingly objective search functions are anything but, and have offered varied approaches to analyzing the systems we so readily adopt. Many have argued for more broad definitions of algorithms to account for the entire socio-cultural process that produces them [Kitchin 2017]. Accordingly, Jamie “Skye” Bianco warns that “tools don’t reflect upon their own making, use or circulation or upon the constraints under which their constitution becomes legible” [Bianco 2012, 99]. Feminist and anti-racism scholars have convincingly shown how algorithmic shortcomings in search, database construction, and knowledge organization can be particularly detrimental to non-mainstream fields. Tara McPherson directly ties the marginalization of race-related studies to “the very designs of our technological systems” and “post-World War II computational culture” that continues to “partition off consideration of race in our work in the digital humanities” [McPherson 2012, 140]. Moya Z. Bailey more broadly asks about the effect on diverse scholars: “How do those outside the categories white and male navigate this burgeoning disciplinary terrain?”, and Roopika Risam questions the degree to which digital humanities processes as a whole (re)produce centers and peripheries [Bailey 2011] [Risam 2015]. Such intersectional approaches recognize the structural, ideological and political forces that contribute to the creation and promulgation of digital library technologies.

Scholars have questioned the role of the database, specifically, as its own configured media object and unit of inquiry, rather than just a neutral tool [Manovich 1999] [Vesna 2007]. Librarians, especially those with technology expertise or digital interests, have likewise begun investigating bias in academic discovery systems and scholarly databases. With the creation of massive digital corpuses of scholarship and archives, providers have looked for ways to provide enriching metadata. (Often described as data about data, metadata in this context is added information about a document or item, often used for discovery.) Unfortunately, classifying contents is rife for the introduction of biases. Jeffrey Daniels noted in 2015 that an Ex Libris discovery tool returned sexist results: a search on stress in the workplace returned only a Wikipedia article on “Women in the workforce”, implying that women and stress were the same thing [Reidsma 2019, 3–4]. Matthew Reidsma’s recent book, *Masked By Trust: Bias in Library Discovery*, points to the additional demands on library discovery systems to support a variety of intellectual inquiries that may be particularly complex, including “big, challenging, often contentious topics” without objectively correct answers, in contrast to mundane generic searches, such as “nearest gas station” [Reidsma 2019, 68–71]. Reidsma’s earlier work on Proquest, an academic digital document provider, found that problematic search results related to likely already-marginalized or politicized topics, including “women, the LGBT community, race, Islam, and mental illness” [Reidsma2016].

As one of a relatively small subset of women’s historians with long-term engagement in both machine learning and feminist scholarship, I may be relatively well placed to undertake a critique of JSTOR’s possible algorithmic bias [Jockers 2013, 123-124] [Block and Newman 2011] [Newman and Block 2006]. Still, I write from the perspective of an academic teacher and researcher with substantial knowledge in digital humanities, but without training in taxonomy or library and information science. Accordingly, rather than a comprehensive analysis of JSTOR’s search and topic-based systems, I offer a targeted and detailed review of the topics applied to scholarship in my area of expertise so that I can base all quantitative and text-based analyses on my in-depth understanding of each work’s content, arguments, and foci. This exploration into JSTOR’s topical indexing system for women’s and race histories points to the problematic ways that technology can misconstrue and marginalize scholarship and suggests areas for needed improvement.

## Finding the Algorithm(s)

JSTOR began as an online “Journal Storage” database, and is today a broader not-for-profit digital library built for academic research. JSTOR touts its availability in 10,215 institutions and 176 countries where it provides access to more than twelve million pieces of academic writing in seventy-five disciplines. It is an indispensable resource to the academic community in U.S. and women’s history.

8

Various online guides to JSTOR describe their topic labeling algorithms in general terms. The JSTOR Thesaurus, apparently launched sometime between 2013 and 2017, provides the controlled vocabulary (standardized terminology) that makes up topics [JSTOR 2017] [JSTOR 2019b]. In 2018, Jabin White, a vice president at ITHAKA, the digital technology not-for-profit that produces JSTOR, described JSTOR’s the creation of the Thesaurus as a way to address the need for descriptive and semantic metadata that now provides “additional value” for libraries and users [White 2018]. The Thesaurus is constructed of seventeen public and corporate-produced vocabularies, and is not available to users online. [JSTOR 2019b]. JSTOR reports that it relied on the software company, Access Innovations, to create the Thesaurus. Access Innovations touts its four decades of taxonomy development experience which allows it to create classifications for customers to “fit both your content and the way your users interact with that content” by working “closely with your subject matter experts”.

9

It is not transparent precisely how JSTOR Thesaurus terms become topics for specific pieces of scholarship. A JSTOR support page explains some details of topic indexing its database: “If a term is present at least three times, it is recognized by the thesaurus and triggers the application of a topic” with “up to 10 topics assigned to” an article or chapter [JSTOR 2019a]. It is difficult to tell from available descriptions what precise system(s) JSTOR is using to create its listed topics. A JSTOR taxonomer explained that it relies on “both auto & manual rule creation”, and specified that they use Wikipedia or DBPedia (which includes over five million entities of structured content from Wikipedia) to provide information for and descriptions of topics. <sup>[1]</sup> These topics do not seem to be produced by probabilistic topic modeling, even though JSTOR Labs’ Text Analyzer uses LDA (Latent Dirichlet Allocation), a popular topic model [Snyder 2012]. Such statistically based models are part of the larger field of probabilistic modeling and automatically learn a set of topics that comprehensively describe a set of documents. In the past decade-plus, topic modeling has been increasingly applied to humanities research questions and texts to find themes in such large corpora without *a priori* subject categories [Block 2006] [Meeks and Weingart 2012]. Topic modeling is particularly good at disambiguation (separating different senses of words) and thematically linking words with allied meanings (e.g.: car, automobile, BMW and Ford would all likely be listed in the same topic). Even though information pages on JSTOR’s Thesaurus explicitly “recognize” the need to distinguish among homonyms, in practice, the topic identification system seems to fall short on effective word disambiguation. For example, it lists the topic of “Charity” first for an article by Jessica Millward about a woman named “Charity Folks” [Millward 2013].

10

The end user can only guess at the precise topic-defining and ranking process. For instance, does JSTOR’s topic identification system rely on Part-of-Speech Tagging (identifying whether each word is a noun, verb, etc) to ascertain each word’s or phrase’s role in a given sentence? Or might it rely on Shallow Parsing to “chunk” parts of sentences (nouns, verbs, etc) with less specificity? Do JSTOR’s topic assignment algorithms identify specific text phrases such as Named Entities (proper nouns, such as individual and organization names) and select other multi-word expressions?

11

I have attempted to describe JSTOR’s topic algorithms not because it is every user’s responsibility to understand controlled vocabularies, rule bases, machine indexing, knowledge bases, and the principles and practices of taxonomy construction. But such human and coding details – where topic information comes from, whether indexing is human or machine curated – can create bias in the results. For example, Jessica Parr, a historian with an MS in Archives management, was “surprised” at JSTOR’s use of Wikipedia/DBPedia: “People have been talking about Wikipedia’s considerable flaws with race and gender topics for several years. And these are flaws that haven’t been fixed. Using a tool like DBPedia means your system is full of these racial and gendered biases” [Lapowsky 2015] [Zevallos 2014].

12

Even without becoming expert coders and taxonomers, scholars who use databases can and should ask questions about the ethics of our research tools. We can learn to probe the ethics of databases in the same ways that scholars

13

have spent the past decade productively interrogating the construction of archives [Falzetti 2015] [Fuentes 2016] [Stoler 2010]. Whose stories may be silenced or misrepresented in various classification systems? How can users begin to understand the impact of knowledge systems on our understanding of the scholarship done in particular subject areas? Analyzing the choices of programmers, taxonomers, and other database system contributors makes clear that academic discovery databases are the result of human decisions rather than any imagined algorithmic neutrality. By focusing on the consequences of JSTOR's algorithmic indexing choices for scholarship in women's African American, Native American, and race histories, this essay offers a first step to rethinking the topical indexing of an influential scholarly database.

## JSTOR's Topical Indexing: Erasure and Misinterpretation

Jennifer Morgan's article, "'Some Could Suckle over Their Shoulder': Male Travelers, Female Bodies, and the Gendering of Racial Ideology, 1500-1770", is one of the most downloaded in the *William and Mary Quarterly* [Morgan 1997]. Personally, I have taught it well over a dozen times since its first publication. Morgan's article focuses on European travelers' constructions of African and Native American women's bodies, showing how European print descriptions and images of women's appearance and behavior helped to build the foundations of race-based slavery.

14

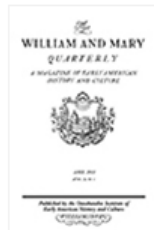


Search JSTOR

All Content



[The William and Mary Quarterly](#) / [Vol. 54, No. 1, Jan., 1997](#) / "Some Could Suckle o..."



JOURNAL ARTICLE

### "Some Could Suckle over Their Shoulder": Male Travelers, Female Bodies, and the Gendering of Racial Ideology, 1500-1770

Jennifer L. Morgan

*The William and Mary Quarterly*

Vol. 54, No. 1, Constructing Race (Jan., 1997), pp. 167-192

Published by: [Omohundro Institute of Early American History and Culture](#)

DOI: 10.2307/2953316

Stable URL: <https://www.jstor.org/stable/2953316>

Page Count: 26

**Topics:** [Men](#), [Black people](#), [African Americans](#), [Children](#), [Academic libraries](#), [Mothers](#), [Childbirth](#), [Civility](#), [Narratives](#)

Figure 1. JSTOR record for Jennifer Morgan, "Some Could Suckle" with Topics.

JSTOR lists nine relevant topics for Morgan's article (Figure 1).<sup>[2]</sup> The first, and presumably most relevant, is "Men". "Women" does not even appear, although "Mothers" does. JSTOR topics further categorize this scholarship as being about "Black people" and "African Americans". Both the present and absent topics for Morgan's article suggest intersectionally biased impacts for historians of slavery, race, African and African American history, and women's history. (On intersectionality, the idea that interlocking systems of power, such as sexism and racism, work together to multiply

15

marginality, see the foundational work of [Crenshaw 1989]). To begin with, the absence of “women” and primacy of “men” is curious, given that JSTOR states that a topic’s “relevance is determined by how frequently the term appears in the piece of content” [JSTOR 2019a].

Given these claims, I undertook an approximate word frequency count of Morgan’s article, wondering if it were possible that the text mentions “men” far more than “women”.<sup>[3]</sup> As Table 1 shows, the occurrence of words such as women/woman/female far outnumber men/man/male. “Women” alone occurs more than seven times as often as “men”. So it is concerning that JSTOR determined that “men” is more relevant a topic than “women”. Moreover, even the category “Black people” misrepresents the article’s focus. A bigram count (occurrences of meaningful two-word compounds) shows that “African women” is the most frequent two-word expression in the article, appearing 32 times, and “black women” is the third most frequent. Calculating bigram frequencies like these offers a methodologically productive way to add specification that better translates broad terms into a topic. Such bigrams show how Morgan’s article focuses on women, not generic people of African descent. But this intersectional identity is erased by JSTOR’s topics. Other researchers have found that Ex libris tools also tended to turn searches related to Africa “into topics about African-Americans”. This may relate to their shared use of Wikipedia for subject information and metadata [Reidsma 2019, 129]. Replication of popular views on race rather than discipline-specific or theoretically informed ones may lead to these kinds of biased results.

16

Women-related terms	Frequency	Men-related terms	Frequency
women	169	men	23
woman	50	man	10
female	30	male	9
TOTAL	241		41

**Table 1.** Select women- and men-related term approximate frequencies in Morgan, “Some Could Suckle.”

The issue with JSTOR’s topical representation of Morgan’s article is not just the absence of women, it is the minimization of an array of terms related to feminist analysis. Versions of the word “mother”, which is sixth on JSTOR’s topic list, appears about 22 times in this article. Yet lexemes of “gender” appear more than 30 times, and “sex” related terms (sexuality, sexual, sexualized) appears more than 50 times, but neither appear as relevant topics. Sexuality and gender are key analytic categories to scholars who do women’s history, making their absence is a notable failing of JSTOR’s topic categorization system.

17

The topics related to the analytic interrogation of racial ideologies and representation of non-white historical actors show additional problems. One of the most relevant analytic terms for historians, “race”, does not feature as a topic, even though the article is about the construction of early modern racial ideologies – as the title clearly conveys. Indeed, the word “race” appears more than twenty times and other forms of the word (racial, racism, racialized) occur almost as often.

18

Even more concerning is the use of “African Americans” as a topic. This is an article primarily about descriptions of African women by Europeans traveling through Africa and Indigenous women in what would become America. In fact, “Africa” appears almost 3 dozen times, and “African” over seventy times, but neither made JSTOR’s topic list. In contrast, “African-American” appears once — in footnote 36 as part of the title of a cited book. Yet JSTOR lists it listed as the third most relevant topic of this article. An entire continent of people has been erased through topical mislabeling in an echo of the Euro-centric bias long critiqued in other library information systems [Howard and Knowlton 2018]. The algorithmic erasure of African and African American women has been repeatedly noted as problematic by scholars across fields [Noble 2018] [Johnson 2018].

19

Likewise, terms frequently related to Native Americans (Native, Indian, Amerindian, Indigenous) occur more than three dozen times in the text, yet did not rate a topic, while “civility”, which appeared just over a dozen times, did. Such topic choices raise the question of whose perspective JSTOR privileges with its topics. Morgan certainly discusses “civility”,

20

but does so in terms of the ways that Europeans mobilized it as a weapon of racemaking. Notions of civility are not a helpful representation of the article's contents because ideas connected to racial ideologies apparently did not merit a JSTOR topic. Markers of civilization have historically marked non-Europeans as exploitable heathens, but the modern meaning of civility as formal politeness elides these racist and colonialist overtones. The erasure of racial ideologies, as well as topics of Africans and Indigenous Americans, means that the necessary topical context here is lost, fundamentally misrepresenting civility's meaning in Morgan's work.

These comparative word frequencies suggest some human-created problems with JSTOR's topics construction. Contrary to what seems to be JSTOR's explanations of its algorithmic processes, these topics are not based simply on word frequency. It appears that JSTOR or outsourced staff have made decisions about what should and should not be in its Thesaurus and the granularity into which some topics should be divided. Unfortunately, those decisions seem to have effects that are both unintended and unattended to. JSTOR's rule base system (which may involve human-curated sets of rules and/or rule based machine learning systems that decide the parameters for classification based on applied domain knowledge) appears to be one that minimizes women, Africans, and scholarship on race as relevant topics. Analytic categorizations that seem most appropriate for scholarship on gender, race, and sexuality, as well as intersectional topics, seem largely absent. This appears to suggest an indexing bias, offering an example "where inaccuracy crosses the line into bias" [Reidsma2016].

21

## Where are the Women in Women's History?

Other women's history articles show similar erasures and misrepresentations. In September 2017, historian Monica Mercado tweeted about the JSTOR topics applied to Linda Kerber's well known article, "Separate Spheres, Female Worlds, Woman's Place: The Rhetoric of Women's History" [Kerber 1988]. Especially since the words "women", "woman", and "female" all appear in the title, Mercado found it surprising that JSTOR's most relevant topic was "Men". A JSTOR's taxonomy manager's response was, effectively, that it was just the algorithm: the result was "relate[d] to how many times those topics appear in the document". She offered as proof that "'Men' appears 63x; Women 25" (Figure 2). The user cannot know exactly what algorithm created those topic frequencies from this brief interaction: did men appear twice as often as women as "topics" or as word frequencies? And "related to" suggests another mediating factors such as human curation or a pre-existing Thesaurus of terms. Regardless, word frequency should have some direct relationship to topic development. Certainly anyone who knows Kerber's work would wonder at JSTOR's topical claim: does one of the founders of U.S. Women's history really focus on men more than twice as often as women in her scholarship?

22



**Sharon Garewal** @sharon\_garewal [Follow](#)

Replying to @monicalmercado

Hi, the 4 topic cards you see highlighted relate to how many times those topics appear in the document, so "Men" appears 63x; Women 25 😊

7:16 AM - 11 Sep 2017

**Figure 2.** Monica Mercado 2017 tweet showing JSTOR's top four topics for Kerber, "Separate Spheres, Female Worlds", article with response from JSTOR Taxonomy Manager.

While I cannot recreate JSTOR's topic-producing algorithms, I can count the approximate frequency of male/female-associated words in Kerber's article as a supplement to my own understanding of its women-focused content. Table 2 shows that "women" appears more than five times as often as "men". In fact, "women" is by far the most frequent word in the article (after stopword removal) with more than 360 mentions. In contrast, "men" appears fewer than 60 times. And an array of women-related words (feminine, feminism, feminist) appear far more than three dozen times. There are zero appearances of any masculine parallel terms (Table 2). As anyone who has read Kerber's work would attest, it is nonsensical that the first topic for this article is "Men".

Term	Frequency	Name	Frequency
women	361	men	55
woman	47	man	10
womanhood	12	manhood	0
female	28	male	23
femine/ism/ist/	42	masculine//ism/ist	0
TOTAL	490		88

**Table 2.** Approximate frequency of select women- and men-related terms in Linda Kerber, "Separate Spheres, Female Worlds, Woman's Place."

Sometime between September 2017 and April 2018 JSTOR attempted to address Mercado's concern over the minimization of women's history. Unfortunately, it appeared to do so by removing "Women" as a topic for this article (See Figure 3). It did add "Women's history" as the fifth topic – but given that this article is entirely about the state of women's history, that seems a rather substantial underrating of its importance, particularly since "Men" remained the first topic.

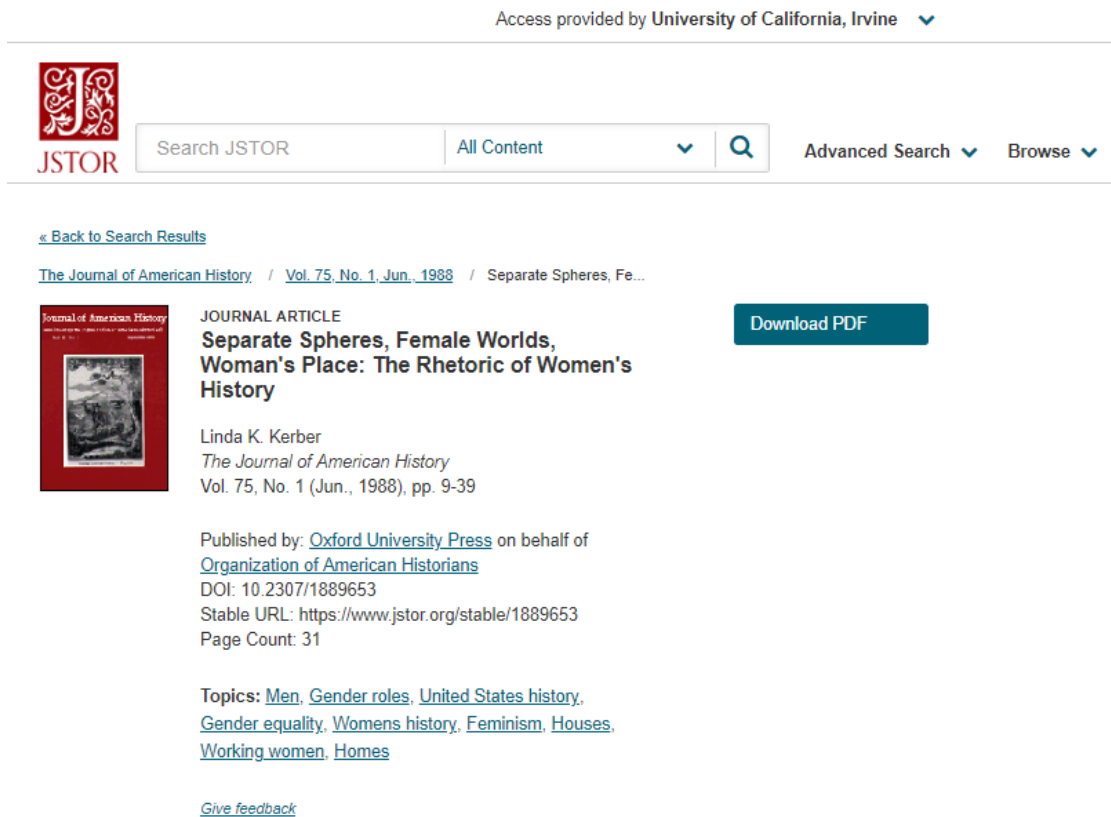


Figure 3. JSTOR categories for Kerber, "Separate Spheres" in June 2018.

When again asked about the absence of the topic of "Women", a JSTOR taxonomist responded on social media that "We removed Women as a topic due to noise a few years ago". (Perhaps she was confused about dates, since Mercado posted that image in September 2017.) The JSTOR staff member was referring to the computer science meaning of "noise" as data that is meaningless or unable to be correctly interpreted. But what does it signal that JSTOR decided that "Women" is "noise", but "Men" is not? This seeming lack of understanding of the historical place of women and women-related topics in academic scholarship suggests that structural power relations – a central analytic accomplishment of the field — are not on JSTOR's radar. One might argue that since "Men" is the standard (or what some call the "null gender"), it might be a less frequent search than "Women", who still tend to carry non-normative status [Wagner et al. 2015]. When men are the subject in the vast majority of historical scholarship, how is it useful for the topic of "Men" to appear as such a relevant topic? And why would "Men" not appear as a major topic for the majority of historical scholarship, then? The same tweet claimed that JSTOR will "probably do the same with Men" (remove it as a too-noisy topic), but as of July 2019, that did not seem to have happened and it remains first in the list of relevant topics for Kerber's article. Men still remains an outsize and inaccurate JSTOR topic in many women's history pieces of scholarship. Regarding women as "noise" effectively erases the hard-won successes of women's history, including the decades of efforts to write women back into historical analysis.

The Kerber article's other assigned topics also minimize the importance of women in the piece. Why would a topic like "US history" be seen as more relevant to this article than "Women's history"? Surely U.S. history is an exceptionally broad, perhaps even a too noisy topical category? Moreover, as any women's historian can attest, "gender" and



“women” are distinct topics of inquiry – this is an article about women far more than gender (a quick frequency comparison: “gender” appears about 2 dozen times, “women” more than 350). Yet “Gender equality” and “Gender roles” both appear as topics. At best, “Gender equality” offers a positivist gloss on Kerber’s piece, which is about understanding women’s lives through the historic construct of separate spheres; not about women achieving gender equality. “Gender roles” is a description of culturally expected behavior. “Gender” as an analytic category is how scholars have theorized the ways that structural sexism allows patriarchy to flourish; in other words, an apparatus to understand gender inequality and oppression. Gender as a category of historical analysis is one of the major inventions of feminist scholarship [Scott 1986]. It denotes far more than women’s roles in a given society. Instead, it is a sophisticated problematization of heteroessential and patriarchal structures of power. Feminist scholars have theorized gender in terms of its performativity, its relation to matrices of domination, and its intersectionality [Butler 2006] [Collins 2008]. Turning gender into “Gender roles” transforms an analytic concept of power relations into a descriptive term that identifies how men and women are expected to behave in a given society.

Other topics show additional inadvertent erasures, suggesting granularity or algorithmic decisions that have substantial consequences for categorization accuracy. The topic of “houses” is puzzling because Kerber’s article does not generally focus on “houses” in the sense of an architectural entity, nor as a woman’s workplace. A frequency count of Kerber’s article shows that “house/houses” appears more than two dozen times — but almost all of those mentions are in reference to “Hull House”, the Chicago settlement house co-founded by Nobel peace prize winner Jane Addams and Ellen Gates Starr. A bigram frequency list confirms that “Hull House” is the fourth most frequent pair of terms. In this case, an algorithmic error has erased the work of a Nobel-prize winning woman rather than offering metadata to promote relevant discovery.

Kerber’s piece is not an exception. In my review of women’s history articles, these kinds of problems appeared repeatedly. For example, Terri Snyder’s 2012 “Refiguring Women in Early American History” is, as its title suggests, a review of the field of early American women’s history [Snyder 2012]. When I first looked at this article’s topics in April 2018, women was its ninth most relevant topic, after the top three of “Native Americans”, “History”, and “African American Literature” (Figure 4). This again raises questions of why “Women” would be seen as noise, but “History” would not – not to mention that it is not accurate to say that Snyder’s article focuses on African American literature.

The social media attention to JSTOR’s shortcomings in April 2018 seems to have led to ad hoc changes. The JSTOR taxonomy manager tweeted that “women” would be “added back as a use case” within the month. (In computational terms, a use case defines the relationship between actors and defined steps; from the user’s perspective, it is how a system will respond to their request. I’m not sure exactly what a “use case” means in this context, since “Women” clearly was a possible JSTOR topic already – just not deemed a highly relevant one.) And indeed, “Women” had moved up to the first, and presumably most relevant topic position by July 2018, and “Women’s history” was now a topic as well, which is a much more appropriate topic than the April topic of “Working women” (Figure 5). Yet we might wonder what domain experts were involved with this decision-making process. It is also worth noting that this shift led to other negative outcomes: all mentions of marginalized groups (African American literature, Native Americans) were removed as topics. In adding a focus on women and gender, JSTOR’s topics eliminated all sense of the article’s intersectional approach to and focus on non-white women in early American history.

ARTICLE

[Refiguring Women in Early American History](#)

Terri L. Snyder

*The William and Mary Quarterly*, Vol. 69, No. 3 (July 2012), pp. 421-450

Download PDF

Add to My Lists

Cite This Item

**Prioritized Terms:** History African American literature

Marriage Native Americans Historiography

**Topics:** [Native Americans](#), [History](#), [African American literature](#), [Historiography](#), [Feminism](#), [American literature](#), [Case histories](#), [Working women](#), [Women](#), [Big history](#).

Figure 4. April 2018 screenshot from JSTOR of Snyder, "Refiguring Women" article topics.



JOURNAL ARTICLE

[Refiguring Women in Early American History](#)

[Terri L. Snyder](#)

*The William and Mary Quarterly*, Vol. 69, No. 3 (July 2012), pp. 421-450

**Topics:** [Women](#), [United States history](#), [Womens history](#), [Men](#), [Auctions](#), [Political revolutions](#), [Gender roles](#), [Womens studies](#), [Swords](#)

Figure 5. July 2018 screenshot from JSTOR of Snyder, "Refiguring Women" article topics. This is a slightly different format from figure 4 due to results list v. individual item view.

While this does suggest that JSTOR investigated and revamped its topic identification in response to critiques, the continued inclusion of "Men" still seems concerning. Moreover, "Swords" and "Auctions" do not seem to be particularly relevant topics to the main arguments of this piece, especially alongside an exceptionally broad topic like "United States history". JSTOR's application of "Sword" as a topic is a recognizable curiosity – and at some point between July 2018 and July 2019, JSTOR apparently recognized its erroneous application; it no longer appears as a topic. However, as with Kerber's article, the addition of "Gender roles" is a less obvious and more damaging mislabeling of a field of study. JSTOR topics have missed capturing crucial theoretical underpinnings and arguments in Snyder's essay, and have seemed to present women's history with rose-colored glasses that not only elide implications of struggle, conflict, and oppression, but in their new formulation, further erase Indigenous and African American women. While JSTOR's responsive efforts are commendable, feminist and social justice workers have long argued that impact matters far more than intent [Utt 2013]. Good intentions still lead to biased algorithmic results when programmers and chosen domain experts do not effectively or appropriately analyze scholarship.

30

The disconnect between JSTOR's topics and what field experts would see as the significant content of these publications reflects ongoing discussions in digital humanities regarding algorithmic mediation and the role of classification versus content representation. With the rise of full text data mining capabilities, Library of Congress subject headings and similar controlled indexing vocabularies may seem to be too broad-brushed an approach for users accustomed to searching the internet and for exact strings of text. Full-text-based topic indexing holds the promise of using an author's words to generate precise subject categorizations rather than slotting new work into *a priori* taxonomies. But as these examples show, more technological mediation does not necessarily lead to better outcomes. JSTOR has seemed to tinker with ways to improve its topic indexing, but it continues to fall short. Women-focused

31

histories are still being categorized as focusing on men. Without topic terms that can capture sophisticated analytic content analysis of race and gender, JSTOR topics continue to misrepresent scholarly content.

## Book Topics: Additional Text, Added Bias

In recent years, JSTOR has expanded its corpus beyond journal articles to include digital versions of monograph and anthology chapters from a variety of academic publishers. This means that any inadvertent sexism and racism in algorithmic topic systems have potentially expanded to pollute book-length scholarship. One might think that longer texts divided into multiple chapters might ameliorate some of the errors of the article topic categorizations. Unfortunately, it appears that similarly problematic topic indexing has propagated these new genres of digital scholarship, expanding the bias presented to users.

32

Early American historian Ann Little has explained that her biography, *The Many Captivities of Esther Wheelright*, aims to move beyond privileging men's recounting of and relationships in women's lives. Little wrote Esther Wheelright's biography to "tell the stories of the girls and women who loved her, clothed and fed her, educated her, worked and prayed with her, competed with her, and buried her" [Little 2016, 12]. But JSTOR's topic categories do not convey Little's women-centered approach. A word cloud made up of the JSTOR topics listed for seven book chapters from *The Many Captivities* visually represent the topics assigned to the book's chapters, with the size of words correlating to the number of times the topic appeared (Figure 6).

33

As Figure 6 shows, neither "Women" nor "Women's history" appear as a topic for any of the chapters. Yet this is clearly a book focused on a woman. If we turn to the book's full text, the most frequent word, appearing more than 1300 times, is "her". In contrast, "his" appears only about 300 times. In fact, of the top-10 most frequently used words in the book manuscript, six are related to women (her, Esther, she, Ursuline (an order of nuns), women, mother/s) and none to men. "Mothers" was identified as the most frequent JSTOR topic for Little's book. Unfortunately, its topic frequency does not disambiguate Little's very different usages of the word: while early sections of the book talk about familial mothering, most of the mentions of "mothers" (more than 300 of the c. 445 occurrences) refers to the head of a female religious community. Indeed, the second most common bigram in the book was "Mother Esther", and similarly, "Sister Marie" was a top-ten most frequent bigram. Ideally, any meaningful topical categorization system would disambiguate word sense to avoid these kinds of pitfalls and omissions. This again may suggest that JSTOR has not effectively evaluated the need for disambiguation to accurately represent complex topics.

34



Figure 6. Wordcloud of all JSTOR topic phrases in Little's *Esther Wheelwright* book chapters.

As in its topical assignments to articles, JSTOR's miscategorizations seem to be particularly problematic in relation to traditionally marginalized groups. Little's book shows this erasure of Indigenous people, specifically. A substantial portion of Little's book focuses on Wheelwright's interactions with the Wabanaki. The word "Wabanaki" appears over 400 times in the text, and "Indian" and "Native" add another 280+ appearances. (For comparison, "Governor" appears fewer than 100 times, but still appears as a JSTOR topic for multiple chapters.) Yet the only ethnicity-related topic JSTOR offers is "White people", which is a topic assigned to three different chapters. The book's chapter that focuses on Wheelwright's captivity in a Wabanaki community where she was known as Mali includes extensive discussion of Wabanaki gender, social, and cultural practices. Yet of the eight topics assigned to that chapter, only one, "Wigwams", relates directly to Indigenous people, even though it is mentioned fewer than two dozen times. This effectively reduces the central role of the Wabanaki and other Indigenous groups to mention of the material object of their housing. Most of the other topics assigned to that chapter relate to Catholicism, which seems to promote a Euro-centric and settler-colonial bias that erases Indigenous people.

These problems harm authors who have spent years thoughtfully framing their scholarship and reasonably expect that databases will allow others to accurately discover their content. When I shared JSTOR's topic results and my word frequency analysis with Ann Little, she responded with deep frustration "that I apparently (successfully!) wrote a biography centering on girls and women's lives — as your word count shows — but all of that effort gets swallowed up by JSTOR's deeply flawed algorithm. How could it so deeply distort my book and misapprehend its purpose?" Moreover, she writes that "the near-total erasure of Wabanaki (Native American, First Nations), French Canadian, and Anglo-American people is also deeply concerning — it's as though JSTOR has its own view of what history is".<sup>[4]</sup> JSTOR's seeming tendency to apply topics that misrepresent content does not appear random. Without fully understanding

JSTOR's topical identification processes, I can only guess that, as Little suggests, it seems to be rooted in a lack of expert knowledge in the fields that it is seeking to topically identify. It may be that JSTOR, which covers an array of disciplines, does not have a system that engages experts from each field of study in the creation of a controlled vocabulary. But the result may be topic terms that misrepresent and misconstrue the content of publications. In this case, the focus of JSTOR's topics on a seemingly monolithic presentation of Eurocentric terms in Little's scholarship erases the intersectional and diverse communities her work presents.

Other books' JSTOR topics show similar shortcomings surrounding content on gender and race. Jennifer Morgan's *Laboring Women: Reproduction and Gender in New World Slavery*, is, as its title suggests, a study of enslaved women in the 16th-18th centuries. JSTOR identified "Men" as a chapter topic word seven times in but "Women" only five times – as is visually striking in Figure 7. "Son" and "Children" are topics, but "Daughter" is not. Topics do include some quasi-conceptual themes, such as the terms "Women's rights" and "Gender equality", but it is hard to see how these are relevant terms for a study of enslaved women. It is also notable that JSTOR does not apply the topics of "Race" or "Racism" to *Laboring Women*. Moreover, the system of "Slavery" appears as a topic only once. Instead, "Slaves" and "Slave ownership" are relatively frequent topics. But these two are not parallel categories. It would make more sense to have a topic of "Slave owners" or "Enslavers" — not the passive construction of "slave ownership" — in opposition to "Enslaved People" or "Slaves". Having such unmatched categories sidesteps the reality that white people owned, traded, and tortured human beings under that "ownership" category and erases the abhorrent power abuses inherent in enslavement. It also suggests a lack of awareness of the current state of the field. Historians of slavery have fought to recognize enslaved people as individuals first, and enslavement as a condition by referring to "enslaved people" rather than "slaves" [Foreman et al.].



Figure 7. Wordcloud of all JSTOR topic phrases in Morgan's *Laboring Women* book chapters.

The absence of conceptual categories such as race/racism raises serious concerns about the classification of historical scholarship. Historians do not just describe people; we make arguments about systems of power. For example, the publisher describes Daniel Livesay's recent book, *Children of Uncertain Fortune: Mixed-Race Jamaicans in Britain and the Atlantic Family* as focusing on "the largely forgotten eighteenth-century migration of elite mixed-race individuals from Jamaica to Great Britain, *Children of Uncertain Fortune* reinterprets the evolution of British racial ideologies as a matter of negotiating family membership" [UNC Press 2018]. This is a book about mixed-race individuals, racial ideologies, and the role of familial relationships across the Caribbean and Great Britain. Its assigned Library of Congress subjects also make clear that the book focuses on race: all 10 LOC subjects assigned to the book include the word "race" or "racially mixed" (Figure 8). Again, this is not to argue that LOC classification schemas are free of bias [See, for example, Berman 2014, Dudley 2015, Hathcoc 2016, Knowlton 2005]. But Library of Congress subjects do provide another support for the centrality of race in this book.



LIBRARY OF CONGRESS  
ONLINE CATALOG

LC Online Catalog Quick Search

#### LC Subjects

Racially mixed people--Jamaica--Social conditions--History--18th century.  
 Racially mixed people--Jamaica--Social conditions--History--19th century.  
 Racially mixed people--Great Britain--Social conditions--History--18th century.  
 Racially mixed people--Great Britain--Social conditions--History--19th century.  
 Racially mixed people--Civil rights--Jamaica--History--18th century.  
 Racially mixed people--Civil rights--Great Britain--History--18th century.  
 Racially mixed people--Civil rights--Jamaica--History--19th century.  
 Racially mixed people--Civil rights--Great Britain--History--19th century.  
 Great Britain--Race relations--History.  
 Jamaica--Race relations--History.

**Figure 8.** Library of Congress Subject Headings for Livesay, *Children of Uncertain Fortune* clearly show the centrality of race and racial mixture to the book's content.

Unfortunately, JSTOR's topics for *Children of Uncertain Fortunes* are problematic. Not only are race, racism, racial mixture, or any related terms completely absent from JSTOR's chapter topics, the most frequent topic phrase associated with the book's contents is "White People" (Figure 9). Moreover, even though much of this book focuses on free people of color, there is no parallel focus on "Black People" or other appropriate terms to identify people of African descent. Instead, non-white people are presumably represented with the JSTOR's repeated use of the topic of "Slaves".



**Figure 9.** Wordcloud of all JSTOR topic phrases in Livesay's *Uncertain Fortune* book chapters.

JSTOR does apply the topic of “African American” once – but this is a book about Jamaica and Britain, not North America, and I could only find “African American” in the text one time (Afro-Caribbean appears somewhat more frequently). According to searches undertaken in the Kindle version of the book, “Black” appears approximately 380 times, and “Slave” about 450 times, but “Black People” did not rate a topic while “Slaves” rated one for multiple chapters. Africa/African appears in the text over 300 times, but also did not rate a category, while “Bequests”, which appeared under 100 times, did. Neither was there any topic related to mixed-race people, despite the term “mixed” (as in mixed-race, mixed heritage, mixed ancestry, etc) appearing almost 800 times. Of course, it is hard to make conclusive arguments about an entire book’s content from simple word counts. But between the book’s description, the LOC subject headings, and the above word frequencies, it does seem that JSTOR has flattened important historical differences into racial binaries. As importantly, it seems to employ algorithms that privilege whiteness, and can only see non-white people in terms of their enslavement, rather than as multidimensional human beings with specific ethnic and racial histories. By not focusing on conceptual categories that are central to historical scholarship, JSTOR’s topics do not effectively allow for discovery related to the central analytic accomplishments of this scholarship.

## Conclusion: Interactions, Reactions, and Ways Forward

JSTOR has explained its “Topics” as experimental and welcomes feedback. Every document’s topic list is accompanied by a thumbs-up/thumbs-down clickable icon and a link to “Let us know!” if users see something inaccurate (Figure 10).





BOOK CHAPTER  
CHAPTER 1 INHERITANCE, FAMILY, AND  
MIXED-RACE JAMAICANS, 1700–1761

Download PDF

pp. 20-89

From the Book

[Children of Uncertain Fortune: Mixed-Race Jamaicans in Britain and the Atlantic Family, 1733-1833](#)

Series: [Published for the Omohundro Institute of Early American History and Culture, Williamsburg, Virginia](#)

Copyright Date: 2018

Published by: [Omohundro Institute of Early American History and Culture, University of North Carolina Press](#)

Page Count: 70

Stable URL:

[http://www.jstor.org/stable/10.5149/9781469634449\\_livesay.8](http://www.jstor.org/stable/10.5149/9781469634449_livesay.8)

Topics: [White people](#), [Marriage](#), [Chambers of commerce](#), [Children](#), [Plantations](#), [Women](#), [Kinship](#), [Bequests](#), [Slaves](#)

Were these topics helpful?   
[See something inaccurate? Let us know!](#)

Figure 10. Example of JSTOR request for feedback on topic lists.

There is no question that JSTOR topic applications have an array of seemingly benign errors. In the scholarship discussed above, we might question the relevance of topics on “Swords” (Figure 4) and “Academic libraries” (Figure 1), for example. One user publicly noted that an article with the phrase “to bear” erroneously was assigned a topic of “Bears”, again suggesting problems with disambiguation of word senses.<sup>[5]</sup> In a response to another user sharing seemingly nonsensical JSTOR tags in March 2018, a JSTOR representative responded that “We reviews [sic] those each weed [sic] to make updates to our rules.”). Besides, I imagine, cursing Twitter’s no-editing-post-tweet rule, JSTOR representatives do seem to take individual error notifications under advisement and seek to improve the system. Matthew Reidsma reported having similar experiences with offering feedback on questionable results to ProQuest [Reidsma2016].

The need to gather user feedback for new categorization systems is understandable. At the same time, what responsibility do organizations – let alone one that charges Universities hundreds of thousands of dollars - have to evaluate even an admittedly experimental system for racism and sexism before widely releasing it [JSTOR 2019c]? JSTOR has known that there are problems surrounding these issues since at least 2017 (See Figure 2). Ad hoc attention to user concerns does not seem to be a best practice. A lacking systemic response to racist or sexist algorithmic results is, unfortunately, all too common [Wachter-Boettcher 2017, 133].

Users clicking through topics should have reasonable expectations of accuracy and lack of bias. While it may be easy to understand that the characterization of “Bears” in a piece discussing bearing arms is incorrect, it is more concerning when topics *seem* accurate, but are actually marginalizing and misrepresenting women and non-white people. JSTOR has created a system that can produce topics, but perhaps has not fully evaluated whether those topics have

42

43

44



appropriate disciplinary meaning, nor whether their use of mainstream sources like Wikipedia or DPedia may have introduced an array of racist and sexist terminology and beliefs.

Moreover, when feminist scholars publicly raised these issues with JSTOR, there seemed to be a sense – in line with their exclamation-pointed “Let us Know!” topic list suggestion – that users will volunteer constructive feedback to help improve their system. JSTOR states it has worked with 30 subject matter experts who “volunteer their assistance” and proclaims that “We are also always happy to talk to Subject Matter Experts about particular vocabularies. If you have suggestions or want to talk to us about the thesaurus, email us” [JSTOR 2019a] [JSTOR 2019b]. Given that users pay (either individually or through an institution) for access to JSTOR’s various databases, this seems particularly problematic. JSTOR may be not-for-profit, but it is not a volunteer organization. And asking female scholars who have already identified issues of sexism and racism in JSTOR’s system to spend more “pro bono” time working on a solution conveys a regrettable lack of understanding of the unrecognized service work regularly expected of women in academia [Guarino and Borden 2017] [Babcock et al. 2018].

45

Safiya Noble astutely advises that we “ask ourselves what it means in practical terms to search for concepts about gender, race, and ethnicity only to find information lacking or misrepresentative, whether in the library database or on the open web” [Noble 2018, 142]. I have no reason to believe that JSTOR’s problems are rooted in purposeful racism and sexism. But ultimately, that is irrelevant to their results. The examples I have offered here of problematic JSTOR topics suggests that its taxonomers and programmers have perhaps not adequately addressed Noble’s questions. If, as it seems, these kind of biases are widespread, JSTOR appears to be abdicating its responsibilities to provide a non-racist and non-sexist product. Such shortcomings when programming a complex system may be understandable but they should not be acceptable. They are structural issues that need to be addressed beyond inviting crowd-source error finding.

46

The topics assigned to the articles and books I analyzed here suggest that JSTOR fails particularly well in reference to marginalized histories: for women, for Africans and African Americans, and for Native Americans. Articles on women’s history are assigned the topic of “Men”, which is not even a particularly relevant topic of analysis in historical study. Scholarship on Africa and people of African descent are miscategorized as being about African Americans, who are then assumed to be relegated to a presumed slave status. Indigenous people are viewed through settler colonial and Eurocentric perspectives. Important conceptual categories like race and gender are elided or erased.

47

Shortcomings in JSTOR’s topic classifications raise an array of ethical questions. JSTOR only exists because scholars’ intellectual labor fills its database. What does JSTOR owe in return when it classifies and categorizes the fruits of that labor? Moreover, JSTOR promotes its topic system as particularly useful for students [JSTOR 2018]. I, as a scholar with decades of domain expertise, can easily look at the categorizations of Jennifer Morgan’s “Some Could Suckle” article and know that it is about Africans, not African Americans. Or that the topic of women is not equivalent to the analytic category of gender. But for students who are an ostensible target audience for topical offerings, JSTOR is providing problematic information.

48

It matters that I came across this problem organically, in the course of reviewing the field of early American women’s history because it suggests that JSTOR’s racist and sexist biases may affect others interested in race and gender as historical constructs. One of the challenges for ever-expanding digital document providers is how to offer useful access to their contents. The staff who are tasked with creating the systems that produce these topics no doubt work to the best of their ability because they believe that knowledge preservation and retrieval methods matter. The solutions, then, are far more complex than “don’t be racist/sexist”. This is not about individual responsibility; it is about structural failings. How we can search relates to the scholarship we can find and the knowledge that we produce.

49

I would not be surprised if the biases I have identified may also reflect broader problems with JSTOR’s topic algorithms, relevant to those outside the fields of early American history. For example, on July 5, 2019, clicking on the JSTOR topic of “Men” returned “Variation in Women’s Success across PhD Programs in Economics” as the most relevant scholarship. Several days later (July 8, 2019), the most relevant article under the topic “Men” was listed as “Women in the medieval wall paintings of Canterbury cathedral”. What does it say about JSTOR’s topic-producing algorithms that

50

scholarship that very much appears to be focused on women are the top results for the category of “Men”? I suspect more research would show other kinds of broad classification problems that shade into bias. For instance, the most relevant result when clicking on the topic “White People” is a chapter from a book on *Martin Luther King Jr’s Theory of Political Service*. The top result of “White American Culture” (a problematic category itself) is an article on “Langston Hughes, African American Literature, and the Religious Futures of Black Studies”. These kinds of miscategorizations risk derailing inexperienced researchers, curtailing the use of JSTOR’s resources, and ultimately making important scholarship less easily discoverable than it should be by categorizing scholarship on women and African Americans as being most relevant to topics about men and white people.

I hesitate to move beyond my critique to offer concrete solutions to JSTOR’s practices both because I have only a limited sense of what JSTOR practices entail, and a much better sense of my own limitations, which include only passing knowledge of classification systems, taxonomy, database design, and the multiple needs of a massive multi-disciplinary project like JSTOR. That said, I can comment as an end user on the ways that JSTOR topics do and do not serve scholars’ research and teaching needs, particularly in light of recent scholarship on such issues in a variety of big data systems.

51

Algorithmic transparency and algorithmic accountability have been burgeoning research areas, as theorists struggle to see the far-reaching consequences of big data machine learning technologies – what some have termed the “social power of algorithms” [Gillespie and Seaver 2015] [Dickey 2017] [Neyland and Möllers 2017]. Some have focused on engineering-type solutions, suggesting ways to improve algorithmic decision making by identifying moments of bias entry (i.e., training data bias, algorithmic focus bias, algorithmic processing bias, transfer content bias, interpretation bias, non-transparency of bias) [Silva and Kenney 2018]. Others have suggested how construction of ethical workflows might ameliorate inadvertently ideological outcomes or a “practical algorithmic ethics” that can be used to analyze the virtues and consequences of individual algorithms [Hepworth and Church 2018] [Sandvig et al. 2016]. While some technology and society scholars have pushed for transparency in computer algorithms, others have noted that transparency does not negate biased outcomes.

52

Some scholars have suggested that wishing for unbiased classification systems is heading down a mistaken path. Feminist and queer studies writers have proposed various theoretical rethinkings of how we view metadata. Librarian Emily Drabinski suggests alternatives to the notion that biases can be corrected and classification systems can be made objective. Instead, she employs queer theory to suggest “new ways of thinking about how to be ethically and politically engaged on behalf of marginal knowledge formation”. She argues that we should to “teach knowledge production as a contested project” so that users recognize and engage the bias in knowledge organization systems, rather than expecting functional solutions to cataloguing bias [Drabinski 2013, 96, 108]. Teaching students numerous search tactics – and how to recognize problematic search results – are valuable cross-disciplinary skills to impart [Grey and Hurko 2012]. While this end-user-interrogation approach is a useful one, it is not likely to be entirely successful, particularly when there is not algorithmic transparency. Drabinski herself recognizes that “privatized corporate algorithms” make information organization “less and less apprehendable” [Drabinski 2016]. Similar to Drabinski’s theoretical approach, Anupam Chander suggests instead a “transparency of inputs and results” that will make visible the discriminatory production, rather than eliminate it [Chander 2017]. Being aware of common classification biases – even if we cannot know their exact production process — offers one path to becoming a more thoughtful user of academic discovery systems.

53

Matthew Reidsma, one of the leading investigators of bias in library discovery systems, agrees with the need for user interrogation, and simultaneously proposes several areas of specific improvements. These include paying attention to the degree to which our searches are powered by proxies for the information we request; less unthinking trust in algorithms; increased diversity among programmers; working toward an algorithmic ethics; and intentional audits of software tools [Reidsma 2019, 148–70]. Similarly, I can imagine practical solutions that reshape the relationship between institutional consumers and database providers. JSTOR, as a not-for-profit organization, may have more obligation than privately owned digital document providers to live up to academic community standards. As the recent University of California resistance to publishing giant, Elsevier, has shown, universities can marshall their considerable power as consumers to insist on a range of standards that are in line with their own values and priorities [McKenzie

54

Just because JSTOR can create topics does not mean that it offers useful metadata across fields of scholarship. It might be valuable to rethink the metadata sources and subject expert review processes that JSTOR currently practices. In the era of PhD training that moves beyond professorial careers, perhaps JSTOR could partner with professional organizations to offer internships that pair PhD students in specific fields with library and information science experts to review subject-based metadata. This might help evaluate the degree to which topics from the JSTOR Thesaurus's controlled vocabulary meet standards for appropriacy (is a given term appropriate for the target audience?) and currency (does it reflect current common usage?). For scholars of women, race and colonialism, at least, the answer currently appears to be no for too many of JSTOR's topic categorizations.

A commitment to changes in and increased transparency of review processes might also be a positive step. JSTOR does make individual changes when users point to errors. But it seems unprepared to deal with a wholesale review of the ways that their topic assignments foreground inadvertent sexism and racism, and may lack structural due diligence against bias. JSTOR prides itself on "enhancing its content with strong metadata" [Humphrey 2019, slide 9]. It is not clear that its current topic system meets that standard. Perhaps involving end users more systematically (beyond thumbs up/thumbs down and ad hoc communications) would promote a more transparent knowledge organization system. Scholars have pointed specifically to the need to "expand the boundaries of LIS [Library Information Systems]" to better understand the "ways in which tools impact the research process" [Manoff 2015, 526]. Ultimately, JSTOR is a crucial resource for historians and other scholars, students, and educators. It may be easy for an outsider to find shortcomings, but finding solutions is far from an individual task. Raising awareness of these kinds of biases can encourage academic communities to work with JSTOR and other digital providers to create systems that better reflect the scholarship on which they build their systems.

## Notes

[1] The JSTOR taxonomer who tweeted in response to questions about Topics has apparently deleted her Twitter account. Screenshots of quoted tweets, as well as of JSTOR search result images not reproduced here are on file with the author.

[2] Unless otherwise specified, the topics assigned to scholarship were taken from JSTOR in June-July 2018. Over time, its topic model results have seemed to change, which I have mentioned when relevant.

[3] Frequency counts included basic text normalization: removal of punctuation; lowercasing of all text; rudimentary de-pluralization; stopword omission; and ignoring of words of fewer than 3 characters. Since all I needed were approximate frequency counts, I did not clean up the text before processing. Thanks to David Newman for technical assistance.

[4] Little, Personal Communication, July 14, 2018.

[5] <https://twitter.com/alforrester/status/927981173049118720>

## Works Cited

- Babcock et al. 2018** Babcock, Linda, Maria P. Recalde, and Lise Vesterlund. "Why Women Volunteer for Tasks That Don't Lead to Promotions". *Harvard Business Review*, July 16, 2018. <https://hbr.org/2018/07/why-women-volunteer-for-tasks-that-dont-lead-to-promotions>.
- Bailey 2011** Bailey, Moya Z. "All the Digital Humanists Are White, All the Nerds Are Men, but Some of Us Are Brave". *Journal of Digital Humanities* 1, no. 1 (Winter 2011). <http://journalofdigitalhumanities.org/1-1/all-the-digital-humanists-are-white-all-the-nerds-are-men-but-some-of-us-are-brave-by-moya-z-bailey/>
- Berman 2014** Berman, Sanford. *Prejudices and Antipathies: A Tract on the LC Subject Heads Concerning People*. Jefferson, N.C.: McFarland and Co., 2014.
- Bethel 1994** Bethel, Kathleen E. "Culture Keepers: Cataloging the Afrocentric Way". *The Reference Librarian*. (July 1994) 21, vol 45-46. 221-240. [https://doi.org/10.1300/J120v21n45\\_21](https://doi.org/10.1300/J120v21n45_21)
- Bianco 2012** Bianco, Jamie "Skye". "This Digital Humanities Which is Not One" in Gold, Matthew K., ed. *Debates in the Digital Humanities*. Minneapolis: University Of Minnesota Press, 2012. 96-112.

- Blei 2012** Blei, David M. "Topic Modeling and Digital Humanities". *Journal of Digital Humanities* no.1 (Winter 2012). <http://journalofdigitalhumanities.org/2-1/topic-modeling-and-digital-humanities-by-david-m-blei/>.
- Block 2006** Block, Sharon. "Doing More with Digitization". *Common-Place* 6, no. 2 (January 2006). <http://commonplace.online/article/doing-more-with-digitization/>.
- Block and Newman 2011** Block, Sharon and David Newman. "What, Where, When, and Sometimes Why: Data Mining Two Decades of Women's History Abstracts". *Journal of Women's History* 23, no. 1 (March 2011): 81–109. <https://doi.org/10.1353/jowh.2011.0001>.
- Bogost 2015** Bogost, Ian. "The Cathedral of Computation". *The Atlantic*, January 15, 2015. <https://www.theatlantic.com/technology/archive/2015/01/the-cathedral-of-computation/384300/>.
- Buolamwini 2016** Buolamwini, Joy. *How I'm Fighting Bias in Algorithms*. November 2016. [https://www.ted.com/talks/joy\\_buolamwini\\_how\\_i\\_m\\_fighting\\_bias\\_in\\_algorithms/up-next](https://www.ted.com/talks/joy_buolamwini_how_i_m_fighting_bias_in_algorithms/up-next).
- Butler 2006** Butler, Judith. *Gender Trouble: Feminism and the Subversion of Identity*. New York: Routledge, 2006.
- Chander 2017** Chander, Anupam. "The Racist Algorithm?." *Michigan Law Review*. 115, no. 6 (2017) 1023-1045. <https://repository.law.umich.edu/mlr/vol115/iss6/13>.
- Christensen 2008** Christensen, Ben. "Minoritization vs. Universalization: Lesbianism and Male Homosexuality in LCSH and LCC". *Knowledge Organization* 35, no. 4 (2008) 229-238. <https://doi.org/10.5771/0943-7444-2008-4-229>.
- Collins 2008** Collins, Patricia Hill. *Black Feminist Thought: Knowledge, Consciousness, and the Politics of Empowerment*. New York: Routledge, 2008.
- Crenshaw 1989** Crenshaw, Kimberle. "Demarginalizing the Intersection of Race and Sex: A Black Feminist Critique of Antidiscrimination Doctrine, Feminist Theory and Antiracist Politics". *University of Chicago Legal Forum* 1 (1989) 139-167. <http://chicagounbound.uchicago.edu/uclf/vol1989/iss1/8>.
- Diao and Cao 2016** Diao, Junli, and Haiyun Cao. "Chronology in Cataloging Chinese Archaeological Reports: An Investigation of Cultural Bias in the Library of Congress Classification". *Cataloging & Classification Quarterly* 54, no. 4 (May 18, 2016): 244–62. <https://doi.org/10.1080/01639374.2016.1150931>.
- Dickey 2017** Dickey, Megan Rose. "Algorithmic Accountability". *TechCrunch* (blog), 2017. <http://social.techcrunch.com/2017/04/30/algorithmic-accountability/>.
- Drabinski 2009** Drabinski, Emily. "Gendered S(h)elves: Body and Identity in the Library". *Women and Environments International*. 78/79 (Fall/Winter 2009) 16-18. [http://www.emilydrabinski.com/wp-content/uploads/2012/06/emily\\_weimag.pdf](http://www.emilydrabinski.com/wp-content/uploads/2012/06/emily_weimag.pdf).
- Drabinski 2013** —. "Queering the Catalog: Queer Theory and the Politics of Correction". *The Library Quarterly* 83, no. 2 (April 2013): 94–111. <https://doi.org/10.1086/669547>.
- Drabinski 2016** —. "WILU 2016". *Emily Drabinski* (blog), June 7, 2016. <http://www.emilydrabinski.com/wilu-2016/>.
- Dudley 2015** Dudley, Michael. "A Library Matter of Genocide, pt. II". *The Decolonized Librarian*. September 30, 2015. <https://web.archive.org/web/20190712212959/https://decolonizedlibrarian.wordpress.com/tag/library-of-congress/>.
- Eubanks 2018** Eubanks, Virginia. *Automating Inequality: How High-Tech Tools Profile, Police, and Punish the Poor*. New York, NY: St. Martin's Press, 2018.
- Falzetti 2015** Falzetti, Ashley Glassburn. "Archival Absence: The Burden of History". *Settler Colonial Studies* 5, no. 2 (April 2015): 128–44. <https://doi.org/10.1080/2201473X.2014.957258>.
- Foreman et al.** Foreman, P. Gabrielle et al. "Writing about 'Slavery'? This Might Help". Accessed July 30, 2018. <https://docs.google.com/document/d/1A4TEdDgYsIX-hlKezLodMIM71My3KTN0zxRv0IQTOQs/mobilebasic>.
- Fuentes 2016** Fuentes, Marisa J. *Dispossessed Lives: Enslaved Women, Violence, and the Archive*. Philadelphia: University of Pennsylvania Press, 2016.
- Gillespie and Seaver 2015** Gillespie, Tarleton and Nick Seaver. "Critical Algorithm Studies: A Reading List". *Social Media Collective* (blog), November 5, 2015. <https://socialmediacollective.org/reading-lists/critical-algorithm-studies/>.
- Granieri 2014** Granieri, Giuseppe. "Algorithmic Culture. 'Culture Now Has Two Audiences: People and Machines.'" *Medium* (blog), April 30, 2014. <https://medium.com/futurists-views/algorithmic-culture-culture-now-has-two-audiences-people-and-machines-2bdaa404f643>.

- Grey and Hurko 2012** Grey, April and Christine R. Hurko. "So You Think You're an Expert: Keyword Searching vs. Controlled Subject Headings". *Codex: the Journal of the Louisiana Chapter of the ACRL* 1, no. 4 (2012) 15–26. <http://journal.acrla.org/index.php/codex/article/view/47/110>.
- Guarino and Borden 2017** Guarino, Cassandra M., and Victor M. H. Borden. "Faculty Service Loads and Gender: Are Women Taking Care of the Academic Family?" *Research in Higher Education* 58, no. 6 (September 2017): 672–94. <https://doi.org/10.1007/s11162-017-9454-2>.
- Hathcock 2016** Hathcock, April. "White Privilege — See Also Library of Congress". At *The Intersection* (blog), November 5, 2016. <https://aprilhathcock.wordpress.com/2016/11/05/white-privilege-the-library-of-congress/>.
- Hepworth and Church 2018** Hepworth, Katherine and Christopher Church. "Racism in the Machine: Visualization Ethics in Digital Humanities Projects". *Digital Humanities Quarterly* 12 no. 4 (2018). <http://www.digitalhumanities.org/dhq/vol/12/4/000408/000408.html>.
- Howard and Knowlton 2018** Howard, Sara A. and Steven A. Knowlton. "Browsing Through Bias: The Library of Congress Classification and Subject Headings for African American Studies and LGBTQIA Studies". *Library Trends* 67, no. 1 (2018): 74–88. doi:10.1353/lib.2018.0026
- Humphrey 2019** Humphrey, Alex. "Enabling New Methods of Discovery - Digital Preservation Virtual Conference - Lawrence Livermore National Laboratory". June 28, 2019. <https://www.slideshare.net/AlexHumphreys1/enabling-new-methods-of-discovery-digital-preservation-virtual-conference-lawrence-livermore-national-laboratory>
- Introna and Nissenbaum 2000** Introna, Lucas D, and Helen Nissenbaum. "Shaping the Web: Why the Politics of Search Engines Matters". *The Information Society* 16 (2000): 169-185. <https://doi.org/10.1080/01972240050133634>
- JSTOR 2017** "The JSTOR Thesaurus: Improving Discovery on the Platform". *JSTOR*. October 12, 2017. <https://about.jstor.org/events/jstor-thesaurus-improving-discovery-platform/>
- JSTOR 2018** "Searching Topic Cards". *JSTOR Support*. June 30, 2018. <https://support.jstor.org/hc/en-us/articles/115005573368-The-JSTOR-Thesaurus-and-Topic-Cards>
- JSTOR 2019a** JSTOR. "The JSTOR Thesaurus and Topic Cards". *JSTOR Support*. Accessed July 6, 2019. <https://support.jstor.org/hc/en-us/articles/115005573368-Searching-Topic-Cards>.
- JSTOR 2019b** JSTOR. "JSTOR's Thesaurus". *JSTOR Guides*. Updated Dec 21, 2017. Accessed July 6, 2019. <https://guides.jstor.org/c.php?g=743025&p=5317792>.
- JSTOR 2019c** "U.S. Universities and Four-Year Colleges". *JSTOR Fees Overview*. Accessed July 6, 2019. <https://www.jstor.org/librarians/fees/us-universities>.
- Jockers 2013** Jockers, Matthew L. *Macroanalysis: Digital Methods and Literary History*. Champaign: University of Illinois Press, 2013.
- Johnson 2018** Johnson, Jessica Marie. "Markup Bodies: Black [Life] Studies and Slavery [Death] Studies at the Digital Crossroads". *Social Text* 36, no. 4 (137) (December 2018): 57–79. <https://doi.org/10.1215/01642472-7145658>.
- Kerber 1988** Kerber, Linda K. "Separate Spheres, Female Worlds, Woman's Place: The Rhetoric of Women's History". *The Journal of American History* 75, no. 1 (1988): 9–39. <https://doi.org/10.2307/1889653>.
- Kitchin 2017** Kitchin, Rob. "Thinking Critically about and Researching Algorithms". *Information, Communication & Society* 20, no. 1 (2017). <https://www.tandfonline.com/doi/abs/10.1080/1369118X.2016.1154087>.
- Knowlton 2005** Knowlton, Steven A. "Three Decades Since *Prejudices and Antipathies*: A Study of Changes in the Library of Congress Subject Headings". *Cataloging & Classification Quarterly* 40, no. 2 (August 2005): 123–45. [https://doi.org/10.1300/J104v40n02\\_08](https://doi.org/10.1300/J104v40n02_08).
- Lapowsky 2015** Lapowsky, Issie. "Meet the Editors Fighting Racism and Sexism on Wikipedia". *Wired*, March 5, 2015. <https://www.wired.com/2015/03/wikipedia-sexism/>.
- Little 2016** Little, Ann M. *The Many Captivities of Esther Wheelwright*. New Haven: Yale University Press, 2016.
- Livesay 2018** Livesay, Daniel. *Children of Uncertain Fortune: Mixed-Race Jamaicans in Britain and the Atlantic Family, 1733-1833*. Chapel Hill: University of North Carolina Press, 2018.
- LoC 2019a** Library of Congress. "Library of Congress Classification". Accessed July 2, 2019. <https://www.loc.gov/catdir/cpsol/lcc.html>.

**LoC 2019b** Library of Congress. "Subject Headings and Genre/Form Terms (Cataloging and Acquisitions at the Library of Congress)". Web page. Accessed July 2, 2019. <https://www.loc.gov/aba/cataloging/subject/>.

**Manoff 2015** Manoff, Marlene. "Human and Machine Entanglement in the Digital Archive: Academic Libraries and Socio-Technical Change". *Portal: Libraries and the Academy* 15, no. 3 (July 2015): 513–30. <https://doi.org/10.1353/pla.2015.0033>.

**Manovich 1999** Manovich, Lev. "Database as Symbolic Form". *Convergence* 5, no. 2 (June 1, 1999): 80–99. <https://doi.org/10.1177/135485659900500206>.

**McKenzie 2019** McKenzie, Lindsay. "University of California Cancels Deal with Elsevier after Months of Negotiations". *Inside Higher Ed* March 1, 2019. <https://www.insidehighered.com/news/2019/03/01/university-california-cancels-deal-elsevier-after-months-negotiations>.

**McPherson 2012** McPherson, Tara. "Why Are the Digital Humanities So White? or Thinking the Histories of Race and Computation". Gold, Matthew K., ed. *Debates in the Digital Humanities*. Minneapolis: University of Minnesota Press, 2012. 189–160.

**Meeks and Weingart 2012** Meeks, Elijah and Scott B. Weingart. "The Digital Humanities Contribution to Topic Modeling". *Journal of Digital Humanities* 2, no. 1 (Winter 2012). <http://journalofdigitalhumanities.org/2-1/dh-contribution-to-topic-modeling/>.

**Millward 2013** Millward, Jessica. "Charity Folks, Lost Royalty, and the Bishop family of Maryland and New York". *The Journal of African American History* 98, no. 1 (2013): 24–47. <https://doi.org/10.5323/jafriamerhist.98.1.0024>.

**Morgan 1997** Morgan, Jennifer L. "'Some Could Suckle over Their Shoulder': Male Travelers, Female Bodies, and the Gendering of Racial Ideology, 1500–1770". *The William and Mary Quarterly* 54, no. 1 (Jan 1997): 167–92. <https://www.jstor.org/stable/2953316?seq=1>.

**NISO 2010** National Information Standards Organization. *Guidelines for the Construction, Format, and Management of Monolingual Controlled Vocabularies*. Approved July 25, 2005. (2010) [https://groups.niso.org/apps/group\\_public/download.php/12591/z39-19-2005r2010.pdf](https://groups.niso.org/apps/group_public/download.php/12591/z39-19-2005r2010.pdf).

**Newman and Block 2006** Newman, David J., and Sharon Block. "Probabilistic Topic Decomposition of an Eighteenth-Century American Newspaper". *Journal of the American Society for Information Science and Technology* 57, no. 6 (April 2006): 753–67. <https://doi.org/10.1002/asi.20342>.

**Neyland and Möllers 2017** Neyland, Daniel, and Norma Möllers. "Algorithmic IF ... THEN Rules and the Conditions and Consequences of Power". *Information, Communication & Society* 20, no. 1 (January 2017): 45–62. <https://doi.org/10.1080/1369118X.2016.1156141>.

**Noble 2018** Noble, Safiya Umoja. *Algorithms of Oppression: How Search Engines Reinforce Racism*. New York: New York University Press, 2018.

**O'Neil 2016** O'Neil, Cathy. *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*. New York: Crown, 2016.

**Pasquale 2015** Pasquale, Frank. *The Black Box Society: The Secret Algorithms that Control Money and Information*. Cambridge: Harvard University Press, 2015.

**Paul 2016** Paul, Kari. "Microsoft Had to Suspend Its AI Chatbot After It Veered Into White Supremacy". *Vice* (blog), March 24, 2016. [https://www.vice.com/en\\_us/article/kb7zdw/microsoft-suspends-ai-chatbot-after-it-veers-into-white-supremacy-tay-and-you](https://www.vice.com/en_us/article/kb7zdw/microsoft-suspends-ai-chatbot-after-it-veers-into-white-supremacy-tay-and-you).

**Peet 2016** Peet, Lisa. "Library of Congress Drops Illegal Alien Subject Heading, Provokes Backlash Legislation". *Library Journal*. June 13, 2016. <https://www.libraryjournal.com?detailStory=library-of-congress-drops-illegal-alien-subject-heading-provokes-backlash-legislation>.

**Reidsma 2019** —. *Masked by Trust: Bias in Library Discovery*. Sacramento, CA: Library Juice Press, 2019.

**Reidsma2016** Reidsma, Matthew. "Algorithmic Bias in Library Discovery Systems". *Zenodo*, March 11, 2016. <https://matthew.reidsrow.com/articles/173>.

**Risam 2015** Risam, Roopika. "Beyond the Margins: Intersectionality and the Digital Humanities". *Digital Humanities Quarterly* 9, no. 2 (September 2015). <http://digitalhumanities.org/dhq/vol/9/2/000208/000208.html>.

**Sandvig et al. 2016** Sandvig, Christian, Kevin Hamilton, Karrie Karahalios and Cedric Langbort. "When the Algorithm Itself Is a Racist: Diagnosing Ethical Harm in the Basic Components of Software". *International Journal of Communication* 10

- (2016), 4972-4990. <https://ijoc.org/index.php/ijoc/article/view/6182/1807>.
- Schwartzapfel 2019** Schwartzapfel, Beth. "Can Racist Algorithms Be Fixed?" *The Marshall Project*, July 1, 2019. <https://www.themarshallproject.org/2019/07/01/can-racist-algorithms-be-fixed>.
- Scott 1986** Scott, Joan W. "Gender: A Useful Category of Historical Analysis". *The American Historical Review* 91, no. 5 (December 1986): 1053–75. <https://doi.org/10.2307/1864376>.
- Silva and Kenney 2018** Silva, Selena and Martin Kenney. "Algorithms, Platforms, and Ethnic Bias: An Integrative Essay". *Phylon* 55, no.1 & 2 (Summer/Winter 2018), 9-37. <https://www.jstor.org/stable/26545017?seq=1>.
- Snyder 2012** Snyder, Terri L. "Refiguring Women in Early American History". *The William and Mary Quarterly* 69, no. 3 (July 2012): 421–50. <http://www.jstor.org/stable/10.5309/willmaryquar.69.3.0421>.
- Snyder 2017** Snyder, Ron. "Under the Hood of the Text Analyzer". *JSTOR Labs Blog*. March 7, 2017. [https://labs.jstor.org/blog/#!/under\\_the\\_hood\\_of\\_text\\_analyzer](https://labs.jstor.org/blog/#!/under_the_hood_of_text_analyzer)
- Stoler 2010** Stoler, Ann Laura. *Along the Archival Grain: Epistemic Anxieties and Colonial Common Sense*. Princeton, NJ: Princeton University Press, 2010.
- UNC Press 2018** Livesay, Daniel. "Children of Uncertain Fortune: Mixed-Race Jamaicans in Britain and the Atlantic Family, 1733-1833". UNC Press Page. <https://uncpress.org/book/9781469634432/children-of-uncertain-fortune/>
- Ulloa 2019** Ulloa, Jennifer. "Algorithms Are Racist. Now What?" *Towards Data Science*, April 14, 2019. <https://towardsdatascience.com/algorithms-are-racist-now-what-53fc130bb203>
- Utt 2013** Utt, Jamie. "Intent vs. Impact: Why Your Intentions Don't Really Matter". *Everyday Feminism*, July 30, 2013. <https://everydayfeminism.com/2013/07/intentions-dont-really-matter/>.
- Vesna 2007** Vesna, Victoria, ed. *Database Aesthetics: Art in the Age of Information Overflow*. Minneapolis: University of Minnesota Press, 2007.
- Wachter-Boettcher 2017** Wachter-Boettcher, Sara. *Technically Wrong: Sexist Apps, Biased Algorithms, and Other Threats of Toxic Tech*. New York: W. W. Norton & Company, 2017.
- Wagner et al. 2015** Wagner, Claudia, David Garcia, Mohsen Jadidi, and Markus Strohmaier. "It's a Man's Wikipedia? Assessing Gender Inequality in an Online Encyclopedia". *International AAAI Conference on Web and Social Media (ICWSM2015)*, Oxford, May 2015. <http://arxiv.org/abs/1501.06307>.
- White 2018** White, Jabin. "JSTOR's Metadata Story". *Metadata 2020*. March 9, 2018. <http://www.metadata2020.org/blog/2018-03-09-jstor-story/>.
- Zevallos 2014** Zevallos, Zuleyka. "Sexism on Wikipedia: Why the #YesAllWomen Edits Matter". *The Other Sociologist* (blog), June 7, 2014. <https://othersociologist.com/2014/06/08/wikipedia-sexism-yesallwomen/>.