

Swarthmore College

## Works

---

English Literature Faculty Works

English Literature

---

12-1-2015

# The Fictionality Of Topic Modeling: Machine Reading Anthony Trollope's Barsetshire Series

Rachel Sagner Burma

*Swarthmore College*, [rburma1@swarthmore.edu](mailto:rburma1@swarthmore.edu)

Follow this and additional works at: <https://works.swarthmore.edu/fac-english-lit>



Part of the [English Language and Literature Commons](#)

[Let us know how access to these works benefits you](#)


---

### Recommended Citation

Rachel Sagner Burma. (2015). "The Fictionality Of Topic Modeling: Machine Reading Anthony Trollope's Barsetshire Series". *Big Data And Society*. Volume 2, Issue 2. DOI: 10.1177/2053951715610591 <https://works.swarthmore.edu/fac-english-lit/286>

This work is brought to you for free and open access by . It has been accepted for inclusion in English Literature Faculty Works by an authorized administrator of Works. For more information, please contact [myworks@swarthmore.edu](mailto:myworks@swarthmore.edu).

# The fictionality of topic modeling: Machine reading Anthony Trollope's Barsetshire series

Big Data & Society  
July–December 2015: 1–6  
© The Author(s) 2015  
DOI: 10.1177/2053951715610591  
bds.sagepub.com  


Rachel Sagner Buurma

## Abstract

This essay describes how using unsupervised topic modeling (specifically the latent Dirichlet allocation topic modeling algorithm in MALLET) on relatively small corpuses can help scholars of literature circumvent the limitations of some existing theories of the novel. Using an example drawn from work on Victorian novelist Anthony Trollope's Barsetshire series, it argues that unsupervised topic modeling's counter-factual and retrospective reconstruction of the topics out of which a given set of novels have been created allows for a denaturalizing and unfamiliar (though crucially not "objective" or "unbiased") view. In other words, topic models are fictions, and scholars of literature should consider reading them as such. Drawing on one aspect of Stephen Ramsay's idea of algorithmic criticism, the essay emphasizes the continuities between "big data" methods and techniques and longer-standing methods of literary study.

## Keywords

Topic modeling, literature, distant reading, literary criticism, novels, machine learning

In the final two paragraphs of *The Last Chronicle of Barset* (1867), the last-published of Anthony Trollope's six-novel series detailing the social lives of country clergymen and the effects of clergymen "on the society of those around them," the authorial narrator says a sad goodbye to the fictional English county of Barsetshire in which all six novels are set:

And now, if the reader will allow me to seize him affectionately by the arm, we will together take our last farewell of Barset and of the towers of Barchester. I may not venture to say to him that, in this country, he and I together have wandered often through the country lanes, and have ridden together over the too-well wooded fields, or have stood together in the cathedral nave listening to the peals of the organ, or have together sat at good men's tables, or have confronted together the angry pride of men who were not good. I may not boast that any beside myself have so realized the place, and the people, and the facts, as to make such reminiscences possible as those which I should attempt to evoke by

an appeal to perfect fellowship. (*The Last Chronicle of Barset*, 2002: 860, 861)

The Trollopian narrator wistfully hopes that, by experiencing all six novels, he and his reader have together constructed a social world more persistent than any world a single novel could create. Published between 1855 and 1867, the six novels are all set in and around the fictional county of Barsetshire and its cathedral town of Barchester. Distinct from one another, separately published, the six novels share a geography, some institutions, and several characters. But relations between the six novels are varied and uneven. Some share a core set of main characters and places (*The Warden* and *Barchester Towers*), while others claim a

Swarthmore College, Swarthmore, PA, USA

## Corresponding author:

Rachel Sagner Buurma, Swarthmore College, LPAC 202, 500 College Ave, Swarthmore, PA 19081, USA.  
Email: rbuurma1@swarthmore.edu



relation to the other novels in the series only through the presence of a few familiar minor characters and a handful of recurring geographic locations (*Doctor Thorne*).<sup>1</sup> Characters who play a central role in one novel tend to reappear in small supporting roles or as merely mentioned names in others.<sup>2</sup> Taken together, the six novels suggest a capacious and shifting social world with blurry boundaries—a social world under construction, perhaps, one constantly in the process of being conjured into being by the combination of author and reader the narrator fantasizes about at the very end of *The Last Chronicle of Barset*. While Trollope wished for readers with a knowledge of the novels “perfect” enough to allow them to “reminisce” along with the authorial narrator about “the place, and the people, and the facts,” I will suggest that the introduction of other kinds of readers—machine and human—who know less rather than more about the history, social world, and formal features of the Barsetshire novels might help critics today read Trollope’s series in a new way.

Literary critics have had a particularly difficult time accounting for small groups of related novels like the Barsetshire series, in part because theories of the novel almost always take the single novel as their main unit of analysis. Theories of the Victorian novel also tend to assume (even when they don’t assert) that novelists like Trollope seek to represent the singular social world of the individual novel as a finished and stable totality. And from this perspective, the Victorian novel’s representation of social totality depends upon its unified, singular, self-enclosed formal totality. Critics imagine this formal totality as secured by a controlling omniscient narrator who sees all, knows all, and describes all from a point above and outside the novel. Such a total coherence clearly can’t be the model for the more partial and contingent connections that join the Barsetshire series into a loose group. And yet critics who do deal with groups of novels, such as advocates of “distant reading”, seek to identify large-scale patterns across hundreds or thousands of novels—a search that tends to end with the discovery of new varieties of structural totality.<sup>3</sup> So while distant reading offers us some new insights into the study of novels, it is equally unsuited to understanding the kinds of middle-distance questions raised by Trollope’s small group.<sup>4</sup>

How might we undo the ingrained habit of reading social totality and formal totality together that both close and distant readings of the novel seem to share? How might we instead find a way of reimagining the forms of the six novels as semi-detached and their social relations as more partial and unfinished? Borrowing technology built for relatively “big” data and turning it on the relatively small 1,396,000-word corpus of the Barsetshire series offers us one path. Running various

iterations of the unsupervised latent Dirichlet allocation topic modeling algorithm in MALLET<sup>5</sup> on their collective 314 chapters generates a number of topics that suggest both expected and unexpected connections between the very different novels in the informal series. And these connections, when tracked back into individual chapters and read by humans rather than machines, offer us (among other things) a look at the Barsetshire novels’ own encoding of the layered histories of the novel’s many attempts to capture social relations and social worlds through testing out different genres.<sup>6</sup>

For example, we can look at the various versions of one topic whose most frequently occurring words are likely to be “letter write read written letters note wrote writing received table paper send answer return judge handed desk pen addressed” (here labeled topic 38) (see Figure 1).<sup>7</sup> Turning to the chapters in which the topic is likely to appear shows that the Barchester series isn’t merely full of letters (See Moody, 2003). It is, of course, but the appearance of these letters, notes, addresses, and envelopes suggests not merely an emphasis on correspondence; it also points to a generic revenant, to the series’ haunting by the ghost of the epistolary novel, or novel-in-letters. One of the most popular novelistic forms during the middle of the 18th century, by the century’s end the epistolary novel had fallen out of favor. By the mid-Victorian moment of the Barchester novels it was a distant—but, as this model helps us see, persistent—memory.

A relatively low-density topic, distributed in drips and drabs throughout the Barchester novels, the “letter write read written letters note” topic thus addresses itself to the past epistolary novel genre trapped inside; we glimpse it in outline, like a bricked-up window in a Victorian renovation of a Georgian house.<sup>8</sup> Read alone, the topic can’t tell us anything about this generic fossil; it suggests only the idea that letter exchange and correspondence is a recurring topic or theme, a part of the novels’ “contents.” But when we examine the “topics in documents” output, we realize that the chapters in which characters exchange letters and worry about unsent notes gesture to that earlier genre and even proffer an alternative configuration for the novel (see Figure 2). The topics in documents output even points to one chapter in which the narrator announces that for the moment he will regress to the genre of the epistolary novel for the length of the chapter.<sup>9</sup>

The generative uncertainty of topic modeling is crucial here, and stems from the enabling assumptions of topic modeling—the counter-factual assumptions upon which the topic modeling algorithm is explicitly and deliberately based. Topics are probabilistically created formations, and the algorithm that generates topic models is based on the enabling—but crucially,



very different model of the social than the kind of formal totalities held out to us by the novel theory we currently possess. As a kind of reader who knows nothing at all about the rich historical, formal, and social contexts within which the Barsestshire novels (like any novels) are embedded, the algorithm offers us a new view—not a more accurate one, but a different one that lets us see and interpret our novels in a denaturalized and different light. Rather than suspending us in a totalizing system or network, it decomposes our novels, taking us backwards into a fictional composition history, towards the other potential, unwritten novels the Barchester series might have been. The algorithm helps us imagine the way any given novel contains within it many unfinished and impossible versions of itself—versions no Victorian author would or could have written. More specifically, in my example, it lets us see how the ghostly epistolary connections that stretch within and between novels in the series could replace or contradict any totalizing vision of the social, any model relying on a formal totality secured by the idea of an omniscient narrator of a single novel. In so doing, it jettisons any finished and final version of the fiction in favor of what we might think of as a kind of counterfactual set of notes. In some sense, we might imagine topics as the notes a (fictional) narrator might have taken towards writing the novel it (or she, or he, or even they) inhabits and over which it so often claims authorial agency.

I've offered a brief and particularly reflexive example of the way a topic model can point us not to the existing "contents" of novels imagined as represented worlds, but rather to the kinds of writing that prepared for or generated the Barsestshire novels. In the context of literary study, I argue, we should train ourselves to read topic models as notes written by nobody rather than "contents" merely poured into fictional form. I want to suggest, that is, that all topics generated from literary corpuses can help take us back to earlier imaginary forms and versions—discarded drafts that authors might have written but didn't, outmoded genres that are fragmentarily recycled within new forms. Topic modeling may be most useful for humanists when we use it this way, as a kind of uncanny, shifting, temporary index to the works we know best, rather than trying to imagine it, as we too often do, only as telling us something about the stable "contents" of large literary corpora. Closely linked to older traditions of indexing literature (from Victorian Bible concordances to Caroline Spurgeon's index to all of Shakespeare's figural language in *Shakespeare's Imagery* to Roberto Busa's *Index Thomisticus*), the algorithm's machinic, non-semantic, probabilistic characteristics can help denaturalize our relationship to literature and our

attachments to the assumptions—about the sociality of literary form, in my example—baked into our favorite theories of the novel.

Not something a human would ever create, a topic model nevertheless perhaps has more in common than we might at first suspect with the probabilistic, counterfactual, human-created fictions we think we know. Although topics can look at first glance like a pre-existing "discourse," that is, what topics generated from novels actually offer us is the ultimate formalist fantasy of the components of the novel's representation of a social world—a set of "topics" that make up the "contents" of a corpus with no leftovers, a nearly perfect correspondence between the materials of the work and the finished work itself.<sup>11</sup> As Stephen Ramsay argues in *Reading Machines*, using algorithms need not propel us towards applying an ersatz scientific and scientific evidentiary standard to literary interpretation, but rather should reveal and perhaps help amplify our already part-algorithmic literary-critical reading practices, the regular sets of protocols and procedures of analog literary criticism with which we are very—perhaps sometimes too—familiar (Ramsay, 2011: 14).<sup>12</sup> It is as fantasies of formalist reading practices, perhaps, that topic models of literary texts can be most helpful to human readers—as denaturalizing indexes or suggestive counterfactual maps that open up new interpretive possibilities.

### Acknowledgements

Valuable conversations about and feedback on the ideas in this paper came from Laura Heffernan, David Mimno, Michael Reay, the members of the Swarthmore College Victorian Novel Research Seminar, particularly Allison Shultes, and the members of the Tri-College Digital Humanities Art of Topic Modeling Seminar.

### Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

### Funding

The author(s) received no financial support for the research, authorship, and/or publication of this article.

### Notes

1. Trollope did not plan the six novels as a formal series from the beginning, but rather belatedly shaped them into a group. Scholars differ on when Trollope, his reviewers, and readers, began to recognize all six as a set; see Poovey (2010).
2. For work towards computational approaches capable of capturing some of the ambiguity and complexity of the

- “formal” and “referential” identity of characters within and across novels, see Bamman et al. (2014). For an overview of the literary-critical significance of this work, see an abstract of the same title at <http://www.ark.cs.cmu.edu/literaryCharacter/> (accessed 2 October 2015).
3. On the distant reading of novels see, for example, Franco Moretti (2000a, 2000b); Style, Inc. (2009); David Elson et al. (2010); Jockers and Mimno (2013); Underwood et al. (2013); Piper and Algee-Hewitt (2014); and Dewitt (2015). Also see Special issue on topic models and the cultural sciences, John Mohr and Petko Bogdanov (eds) <http://dx.doi.org/10.1016/j.poetic.2013.08.005> (accessed 2 October 2015).
  4. For a different argument about novel theory and social totality in the Barsestshire novels, see Poovey (2010).
  5. See: <http://mallet.cs.umass.edu/>
  6. For more on MALLET and LDA topic modeling, see the *Journal of Digital Humanities*’ special issue on topic modeling 2:1 (Winter 2012), especially David Mimno’s “The Details: Training and Validating Big Models on Big Data.” Available at: <http://journalofdigitalhumanities.org/2-1/> (accessed 2 October 2015). See also Underwood (2012).
  7. The results to which I refer in the article were created with the same corpus using latent Dirichlet allocation in the MALLET topic modeling tool. The models I refer to here are not finished results, but early (and, close readers will note, messy) snapshots from part of a larger in-progress project. For reference, here is a set of the Barsestshire novels loaded into an implementation of latent Dirichlet allocation in javascript written by David Mimno: <http://rachelsagnerbuurma.org/Barsestshire/topics.html>. For the stopword list and documents used, see: <https://github.com/rbuurma/BarchesterAssumptions>. To use Mimno’s software yourself, see: <https://github.com/mimno/jsLDA>
  8. For an approach with related goals that uses computation to trace linguistic “topologies” that allow us to visualize “a latent sense of one text’s presence in another” (p. 156), “not to make definitive pronouncements about absolute affinities or positions but to identify relative connections that could be otherwise” (p. 159), see Piper and Algee-Hewitt (2014).
  9. Note also that the single chapter is composed of only .05 percent of this topic—in some visualizations of these topics across a larger corpus it might not appear at all (another argument in favor of topic modeling as a technique for indexing middle-distance sets of texts).
  10. As Boyd-Graber et al. note, “Topic models are based on a generative model that clearly does not match the way humans write. However topic models are often able to learn meaningful and sensible models” (2014: 15).
  11. For a comparison between the technique of keyword search as a way to locate “discourses” across multiple texts and the technique of topic modeling as a way of tracking “discourses” see Underwood (2014: 66).
  12. See also Piper and Algee-Hewitt on how their algorithmically enabled identification of related clusters of pages from a subset of the works of Goethe is “not a replacement of hermeneutic reading but its facilitator—part of the long history of technologically informed reading practices” (2014: 162).

## References

- Bamman D, Underwood T and Smith N (2014) A Bayesian mixed effects model of literary character. In: *Proceedings from the 52nd annual meeting of the association for computational linguistics (ACL 2014)*, Baltimore, Maryland, USA, 23–25 June. Available at: <http://acl2014.org/acl2014/P14-1/pdf/P14-1035.pdf> (accessed 20 March 2015).
- Boyd-Graber J, Mimno D and Newman D (2014) Care and feeding of topic models: Problems, diagnostics, and improvements. In: Airoldi EM, Blei D, Erosheva EA, et al. (eds) *The Handbook of Mixed Membership Models and Their Applications*. Boca Raton, FL: CRC Press, Taylor & Francis Group.
- Dewitt A (2015) Advances in the visualization of data: The network of genre in the Victorian Periodical Press. *Victorian Periodicals Review* 48(2): 161–182.
- Elson D, Dames N and McKeown K (2010) Extracting social networks from literary fiction. In: *Proceedings of the 48th annual meeting of the association for computational linguistics (ACL 2010)*, Uppsala, Sweden, 11–16 July, pp. 138–147. Available at: <https://www.aclweb.org/anthology/P10/P10-1015.pdf> (accessed 2 October 2015).
- Jockers M and Mimno D (2013) Significant themes in 19th-century literature. *Poetics* 41(6): 750–769.
- Moody E (2003) Partly told in letters: Trollope’s story-telling art. In: *Ellen Moody’s Website*. Available at: <http://www.jimandellen.org/trollope/partly.told.in.letters.html> (accessed 20 March 2015).
- Moretti F (2000a) Conjectures on world literature. *New Left Review*. Available at: <http://newleftreview.org/II/1/franco-moretti-conjectures-on-world-literature> (accessed 2 October 2015).
- Moretti F (2000b) The slaughterhouse of literature. *Modern Language Quarterly* 61(1): 207–227.
- Piper A and Algee-Hewitt M (2014) The Werther effect I: Goethe, objecthood, and the handling of knowledge. In: Erlin M and Tatlock L (eds) *Distant Readings: Topologies of German Culture in the Long Nineteenth Century*. Rochester, NY: Camden House, pp. 155–184.
- Poovey M (2010) Trollope’s Barsestshire series. In: Dever C and Niles L (eds) *The Cambridge Companion to Anthony Trollope*. Cambridge: Cambridge University Press, pp. 31–43.
- Ramsay S (2011) *Reading Machines: Towards an Algorithmic Criticism*. Urbana: University of Illinois Press.
- Style, Inc (2009) Reflections on 7000 titles (British Novels, 1740–1850). *Critical Inquiry* 36(1): 134–358.
- Trollope A and Gilmartin S (eds) (2002 [1876]) *The Last Chronicle of Barset*. London: Penguin Classics.
- Underwood T (2012) Topic modeling made just simple enough. In: *The Stone and the Shell*. Available at: <http://tedunderwood.com/2012/04/07/topic-modeling-made-just-simple-enough/> (accessed 20 March 2015).

Underwood T (2014) Theorizing research practices we forgot to theorize twenty years ago. *Representations* 127(1): 64–72.

Underwood T, Black ML, Auvil L, et al. (2013) Mapping mutable genres in structurally complex volumes. In:

*Proceedings of the 2013 IEEE international conference on big data*, Santa Clara, CA, USA, 6–9 October. Available at: <http://arxiv.org/abs/1309.3323> (accessed 20 March 2015).

This article is part of a special theme on *Colloquium: Assumptions of Sociality*. To see a full list of all articles in this special theme, please click here: <http://bds.sagepub.com/content/colloquium-assumptions-sociality>.