# Contents

CHAPTER 1

# Introduction

This is a book about the use of texts and language to make inferences about human behavior. Our framework for using text as data is aimed at a wide variety of audiences—from informing social science research, offering guidance for researchers in the digital humanities, providing solutions to problems in industry, and addressing issues faced in government. This book is relevant to such a wide range of scholars and practitioners because language is an important component of social interaction—it is how laws are recorded, religious beliefs articulated, and historical events reported. Language is also how individuals voice complaints to representatives, organizers appeal to their fellow citizens to join in protest, and advertisers persuade consumers to buy their product. And yet, quantitative social science research has made surprisingly little use of texts—until recently.

Texts were used sparingly because they were cumbersome to work with at scale. It was difficult to acquire documents because there was no clear way to collect and transcribe all the things people had written and said. Even if the texts could be acquired, it was impossibly time consuming to read collections of documents filled with billions of words. And even if the reading were possible, it was often perceived to be an impossible task to organize the texts into relevant categories, or to measure the presence of concepts of interest. Not surprisingly, texts did not play a central role in the evidence base of the social sciences. And when texts were used, the usage was either in small datasets or as the product of massive, well-funded teams of researchers.

Recently, there has been a dramatic change in the cost of analyzing large collections of text. Social scientists, digital humanities scholars, and industry professionals are now routinely making use of document collections. It has become common to see papers that use millions of social media messages, billions of words, and collections of books larger than the world's largest physical libraries. Part of this change has been technological. With the rapid expansion of the internet, texts became much easier to acquire. At the same time, computational power increased—laptop computers could handle computations that previously would require servers. And part of the change was also methodological. A burgeoning literature—first in computer science and computational linguistics, and later in the social sciences and digital humanities—developed tools, models, and software that facilitated the analysis and organization of texts at scale.

Almost all of the applications of large-scale text analysis in the social sciences use algorithms either first developed in computer science or built closely on those developments. For example, numerous papers within political science—including many of our

own—build on topic models (Blei, Ng, and Jordan, 2003; Quinn et al., 2010; Grimmer, 2010; Roberts et al., 2013) or use supervised learning algorithms for document classification (Joachims, 1998; Jones, Wilkerson, and Baumgartner, 2009; Stewart and Zhukov, 2009; Pan and Chen, 2018; Barberá et al., 2021). Social scientists have also made methodological contributions themselves, and in this book we will showcase many of these new models designed to accomplish new types of tasks. Many of these contributions have even flowed from the social sciences to computer science. Statistical models used to analyze roll call votes, such as Item Response Theory models, are now used in several computer science articles (Clinton, Jackman, and Rivers, 2004; Gerrish and Blei, 2011; Nguyen et al., 2015). Social scientists have broadly adapted the tools and techniques of computer scientists to social science questions.

However, the knowledge transfer from computer science and related fields has created confusion in how text as data models are applied, how they are validated, and how their output is interpreted. This confusion emerges because tasks in academic computer science are different than the tasks in social science, the digital humanities, and even parts of industry. While computer scientists are often (but not exclusively!) interested in information retrieval, recommendation systems, and benchmark linguistic tasks, a different community is interested in using "text as data" to learn about previously studied phenomena such as in social science, literature, and history. Despite these differences of purpose, text as data practitioners have tended to reflexively adopt the guidance from the computer science literature when doing their own work. This blind importing of the default methods and practices used to select, evaluate, and validate models from the computer science literature can lead to unintended consequences.

This book will demonstrate how to treat "text as data" for *social science tasks* and *social science problems*. We think this perspective can be useful beyond just the social sciences in the digital humanities, industry, and even mainstream computer science. We organize our argument around the core tasks of social science research: *discovery*, *measurement*, *prediction*, and *causal inference*. Discovery is the process of creating new conceptualizations or ways to organize the world. Measurement is the process where concepts are connected to data, allowing us to describe the prevalence of those concepts in the real world. These measures are then used to make a causal inference about the effect of some intervention or to predict values in the future. These tasks are sometimes related to computer science tasks that define the usual way to organize machine learning books. But as we will see, the usual distinctions made between particular types of algorithms—such as supervised and unsupervised—can obscure the ways these tools are employed to accomplish social science tasks.

Building on our experience developing and applying text as data methods in the social sciences, we emphasize a sequential, iterative, and inductive approach to research. Our experience has been that we learn the most in social science when we refine our concepts and measurements iteratively, improving our own understanding of definitions as we are exposed to new data. We also learn the most when we consider our evidence sequentially, confirming the results of prior work, then testing new hypotheses, and, finally, generating hypotheses for future work. Future studies continue the pattern, confirming the findings from prior studies, testing prior speculations, and generating new hypotheses. At the end of the process, the evidence is aggregated to summarize the results and to clarify what was learned. Importantly, this process doesn't happen within the context of a single article or book, but across a community of collaborators.

This inductive method provides a principled way to approach research that places a strong emphasis on an evolving understanding of the process under study. We call this

understanding theory—explanations of the systematic facets of social process. This is an intentionally broad definition encompassing formal theory, political/sociological theory, and general subject-area expertise. At the core of this book is an argument that scholars can learn a great deal about human behavior from texts but that to do so requires an engagement with the context in which those texts are produced. A deep understanding of the social science context will enable researchers to ask more important and impactful questions, ensure that the measures they extract are valid, and be more attentive to the practical and ethical implications of their work.

We write this book now because the use of text data is at a critical point. As more scholars adopt text as data methods for their research, a guide is essential to explain how text as data work in the social sciences differs from its work in computer science. Without such a guide, researchers outside of computer science solving problems run the risk of applying the wrong algorithms, validating the wrong quantities, and ultimately making inferences not justified by the evidence they have acquired.

We also focus on texts because they are an excellent vehicle for learning about recent advances in machine learning. The argument that we make in this book about how to organize social science research applies beyond texts. Indeed, we view our approach as useful for social science generally, but particularly in any application where researchers are using large-scale data to discover new categories, measure their prevalence, and then to assess their relationships in the world.

## 1.1  How This Book Informs the Social Sciences

A central argument of this book is that the goal of text as data research differs from the goals of computer science work. Fortunately, this difference is not so great that many of the tools and ideas first developed in other fields cannot be applied to text as data problems. It does imply, however, that we have to think more carefully about what we learn from applying those models.

To help us make our case, consider the use of texts by political scientist Amy Catalinac (Catalinac, 2016a)—a path-breaking demonstration of how electoral district structure affects political candidates' behavior. We focus on this book because the texts are used clearly, precisely, and effectively to make a social science point, even though the algorithm used to conduct the analysis comes from a different discipline. And importantly, the method for validation used is distinctively social scientific and thorough.

Catalinac's work begins with a puzzle: why have Japanese politicians allocated so much more attention to national security and foreign policy after 1997, despite significant social, political, and government constraints on the use of military and foreign policy discussions put in place after World War II? Catalinac (2016a) argues that a 1994 reform in how Japanese legislators are elected explains the change because it fundamentally altered the *incentives* that politicians face. Before the 1994 reform, Japanese legislators were elected through a system where each district was represented by multiple candidates and each party would run several candidates in each district trying to get the majority of the seats. Because multiple candidates from the same party couldn't effectively compete with their co-partisans on ideological issues, representatives tried to secure votes by delivering the most pork—spending that has only local impact, such as for building a bridge—to the district as possible. The new post-1994 reform system eliminated multi-member districts and replaced them with a parallel system: single-member districts—where voters cast their ballot for a candidate—and representatives

デフレ克服　経済活性化

治安対策に全力　犯罪対策の強化
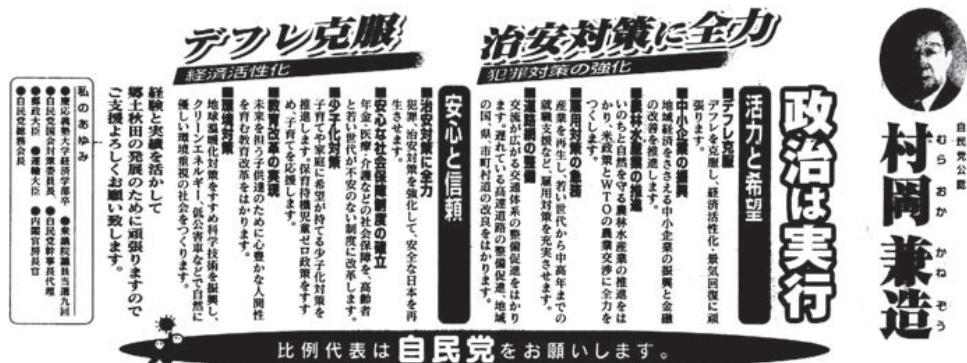
安心と信頼

活力と希望

政治は実行

村岡 兼造

比例代表は自民党をお願いします。

Figure 1.1. An example of a candidate manifesto of Kanezo Muraoka from 2003, Figure 3.7 from Catalinac (2016a).

for the whole country—where voters cast their ballot for a party and the elected officials are chosen from the party's list. This new system allowed the parties to impose stricter ideological discipline on their members and the choices of voters became less about the individual personalities and more about party platforms. Thus, the argument goes, the reform changed the legislators' *incentives*. Focusing on local issues like pork was now less advantageous than focusing on national issues like foreign policy.

The argument proceeds through iteration and induction. To begin understanding the effect of the change in electoral rules on electoral strategy, Catalinac collected an original dataset of 7,497 Japanese Diet candidate manifestos. The manifestos are nearly ideal data for her study: they are important to candidates and voters, under the control of candidates, and available for all candidates for all elections for a period before and after the shift in electoral rules. We discuss the principles for data collection in Chapter 4, but Catalinac's exemplary work shows that working with text data does not mean that we must opt for the most convenient data. Rather, Catalinac engaged in a painstaking data collection process to find the manifestos through archival visits and digitize them through manual transcription. This process alone took years.

With the data in hand, Catalinac uses an inductive approach to learn the categories in her data she needs to investigate her empirical puzzle: what elected officials are discussing when they run for office. Catalinac uses a well-known statistical model, *Latent Dirichlet Allocation* (LDA)—which we return to in Chapter 13—to discover an underlying set of topics and to measure the proportion of each manifesto that belongs to each topic. As Catalinac describes,

> Typically, the model is fit iteratively. The researcher sets some number of topics; runs the model; ascertains the nature of the topics outputted by reading the words and documents identified as having high probabilities of belonging to each of the topics; and decides whether or not those topics are substantively meaningful.… My approach was also iterative and guided by my hypotheses.
> (Catalinac, 2016a, p. 84)

As we describe in Chapter 4, discovery with text data does not mean that we begin with a blank slate. Catalinac's prior work, qualitative interviews, and expertise in Japanese politics helped to shape the discoveries she made in the text. We can bring

this prior knowledge to bear in discovery; theory and hunches play a role in defining our categories, but so too does the data itself.

Catalinac uses the model fit from LDA to measure the prevalence of candidates' discussions of pork, policy, and other categories of interest. To establish which topics capture these categories, Catalinac engages in extensive validation. Importantly, her validations are not the validations most commonly conducted in computer science, where LDA originated. Those validations tend to focus on how LDA functions as a language model—that is, how well it is able to predict unseen words in a document. For Catalinac's purposes, it isn't important that the model can predict unseen words—she has all the words! Instead, her validations are designed to demonstrate that her model has uncovered an organization that is interesting and useful for her particular social scientific task: assessing how a change in the structure of districts affected the behavior of candidates and elected officials. Catalinac engages in two broad kinds of validation. First, she does an in-depth analysis of the particular topics that the model automatically discovers, reading both the high probability words the model assigns to the topic and the manifestos the model indicates are most aligned with each topic. This analysis assures the reader that her labels and interpretations of the computer-discovered topics are both valid and helpful for her social scientific task. Second, she shows that her measures align with well-known facts about Japanese politics. This step ensures that the measures that come from the manifestos are not idiosyncratic or reflecting a wildly different process than that studied in other work. It also provides further evidence that the labels Catalinac assigns to texts are valid reflections of the content of those texts.

Of course, Catalinac is not interested in just categorizing the texts for their own sake—she wants to use the categories assigned to the texts as a source of data to learn about the world. In particular, she wants to estimate the causal effect of the 1994 electoral reform on the shift in issues discussed by candidates when they are running. To do this, she uses her validated model and careful research design to pursue her claim that the electoral reform causes average candidates to shift from a focus on pork to a focus on national security. This is a particularly challenging setting for causal inference, because the reform changes across all districts at the same time. After showing that, in practice, there is a substantial increase in the discussion of national security following the 1994 reforms, Catalinac moves to rule out alternative explanations. She shows that there is no sudden influx of candidates that we would expect to discuss national security. Nor, she argues, does this increase in the importance of national security merely reflect an ideological shift in the parties. And she argues that there is no evidence that voters suddenly want candidates who prioritize national security.

Our brief examination of Catalinac (2016a) reveals how sequence, iteration, and induction can lead to substantively interesting and theoretically important research. Further, Catalinac illustrates a point that we will return to throughout the book, that validations for text as data research are necessary and look quite different from validations in computer science. Rather than a focus on prediction, text as data researchers are much more interested in how well their models provide insights into concepts of interest, how well measurement tools sort documents according to those rules, and how well the assumptions needed for accurate causal inference or prediction are met. These points travel well beyond political science, to other social scientists studying human behavior including sociology (DiMaggio, 2015; Evans and Aceves, 2016; Franzosi, 2004), economics (Gentzkow, Kelly, and Taddy, 2019), psychology (Schwartz et al., 2013), and law (Livermore and Rockmore, 2019).

## 1.2   How This Book Informs the Digital Humanities

Our view of how to apply text as data methods was developed and refined through
our experience with social science research. But we will argue that our approach to text
as data can provide useful insights into other fields as well. In parallel to the meteoric
rise of text as data methods within the social sciences, there has been rapidly grow-
ing interest in using computational tools to study literature, history, and the humanities
more generally. This burgeoning field, termed *Digital Humanities*, shares much in com-
mon with text as data in the social sciences in that it draws on computational tools to
answer classic questions in the field.

The use of text as data methods has drawn considerable funding and has already
made impressive contributions to the study of literature (Jockers, 2013; Piper, 2018;
Underwood, 2019). Computational tools have been used to study the nature of genres
(Rybicki and Eder, 2011), poems (Long and So, 2016), the contours of ideas (Berry and
Fagerjord, 2017), and many other things (Moretti, 2013). To reach their conclusions,
scholars working in this area follow many of the same procedures and use similar tools
to those in the social sciences. They represent their texts using numbers and then apply
models or algorithms that originate in other fields to reach substantive conclusions.

Even though scholars in the Digital Humanities (DH) come from a humanistic tra-
dition, we will show how the goals of their analysis fit well within the framework of
our book. And as a result, our argument about how to use text as data methods to
make valid inferences will cover many of the applications of computational tools in the
humanistic fields. A major difference between DH and the social sciences is that digital
humanists are often interested in inferences about the particular text that is being stud-
ied, rather than the text as an indicator of some other, larger process. As a result, digital
humanities have thus far tended to focus on the discovery and measurement steps of the
research process, while devoting less attention to making causal inferences or predic-
tions. Digital humanists use their large corpora to make new and important discoveries
about organizations in their texts. They then use tools to measure the prevalence of
those quantities, to describe how the prevalence of the characteristics has changed over
time, or to measure how well defined a category is over time.

As with any field that rises so suddenly, there has been considerable dissent about the
prospect of the digital humanities. Some of this dissent lies well outside of the scope of
our book and focuses on the political and epistemological consequences of opening up
the humanities to computational tools. Instead we will engage with other critiques of
digital humanities that stipulate to the "rules" laid out in computational papers. These
critics argue that the digital humanities is not capable of achieving the inferential goals
it lays out and therefore the analysis is doomed from the start. A recent and prominent
objection comes from Da (2019), who summarizes her own argument as,

> In a nutshell the problem with computational literary analysis as it stands is that
> what is robust is obvious (in the empirical sense) and what is not obvious is not
> robust, a situation not easily overcome given the nature of literary data and the
> nature of statistical inquiry.                                    (Da, 2019, p. 45)

Da (2019)'s critique goes to the heart of how results are evaluated and relies heav-
ily on procedures and best practices imported from computer science (as does, it is
worth noting, much of the work she is critiquing). As we have argued above, directly

importing rules from other fields to studying texts in new domains can be suboptimal. When we directly import the recommendations from computer science and statistics to text-based inferences in the humanities or social sciences we might make problematic inferences, recommendations that are misguided, or misplaced assessments about the feasibility of computational analysis for a field.

Yet Da's critique is a useful foil for illuminating a key feature of our approach that departs from much of the work in the digital humanities. In Chapter 2, we offer six core principles which reflect a broader "radically agnostic" view of text as data methods. We reject the idea that models of text should be optimized to recover one true underlying, inherent organization in the texts—because, we argue, no one such organization exists. In much of the digital humanities, and Da's critique, there is an implicit assumption that the statistical models or algorithms are uncovering an ideal categorization of the data that exists outside of the research question asked and the models estimated. This approach is in tension with much of the theoretical work in the humanities, but seemingly arises because this is a motivating assumption in much of computer science and statistics, where it provides a convenient fiction for evaluating model performance.

On our account, organizations are useful if they help us to uncover a categorization of the data that is useful for answering a research question. If two models disagree on how to categorize texts, there is no sense in determining which one is any more "right" than the other. We would not, for example, want to argue that an organization of texts based on the expression of positive or negative emotion is more right than an organization based on the topic of the text. Rather, we will argue that some organizations are more useful than others for addressing a particular question. For example, we might argue that a model is particularly useful for studying genre, because it provides an organization that leads the researcher to an insight about the trajectory of books that would have been impossible otherwise. Once you have an organization, you can find the best *measurement* of that particular categorization. You can then test the measurement with extensive *validation*. But because there is a multiplicity of useful and valid organizations, a method that does not provide a "robust" answer to how texts should be organized will be less concerning than critics argue. What becomes important is the credibility of the validations once an organization has been selected and its utility in answering the research question.

We also will emphasize throughout our book that text as data methods should not displace the careful and thoughtful humanist. And there is no sense in which inferences should be made in the field of digital humanities without the reader directly involved. This emphasis on using computational methods to improve inferences will help allay some concerns about the role of digital humanities scholarship. The computational tools should not replace traditional modes of scholarship. When used well, computational tools should help provide broader context for scholars, illuminate patterns that are otherwise impossible to identify manually, and generally amplify—rather than replace—the human efforts of the scholars using them.

## 1.3  How This Book Informs Data Science in Industry and Government

Computational tools have also revolutionized how companies use text as data in their products and how government uses text to represent the views of constituents. The applications of these tools are nearly endless in industry. Companies use messages that users post on their website to better target advertisements, to make suggestions about

new content, or to help individuals connect with elected officials. In government, there is the chance to use text as data methods to better represent the views of constituents publicly commenting on proposed rule changes at bureaucratic agencies or expressing their views to elected officials.

The stakes are high when applying text as data methods to industrial-scale problems. Perhaps the most politically sensitive application of text as data methods is content moderation: the attempt by social media companies (and sometimes governments) to regulate the content that is distributed on their platform. In the wake of the Russian misinformation campaign in the 2016 US election, social media companies faced increased pressure to identify and remove misinformation from their sites, to report on the effect of misinformation that occurred during the campaign, and to demonstrate that new procedures were fair and did not disproportionately target particular ideologies. The tools used to identify this content will appear throughout this book and will draw on a similar set of computational resources that we introduce.

Beyond the questions of political sensitivity, the application of text as data methods will also be high stakes because of the large amounts of money that will be spent based on the recommendations of the systems. For example, trading firms now use computational tools to guide their investments or to quickly learn about content from central bankers. Text as data methods also help drive advertising decisions that represent a massive share of the economy. Getting these decisions "right," then, is important for many business practices.

Our book is useful for data scientists, because these tasks are inherently social science tasks. Moderating content to suppress misinformation or hate speech is fundamentally a measurement task. When companies decide which ads will cause the largest increase in sales for their clients, they are engaged in causal inference. And when traders make decisions based on the content of documents or statements from officials, they are engaged in prediction. Recognizing the omnipresence of social science within industry is essential, because many data scientists receive their professional training outside of the social sciences. These fields do an excellent job of providing the computational tools necessary for working with the massive datasets that companies create, but often fail to expose researchers to core design principles behind the tasks those tools are built for.

This book, and indeed its very organizational structure, is designed to remove focus from the individual models and computational tools and refocus on the differences between tasks like discovery and measurement or prediction and causal inference. Identifying these differences is essential, because the different tasks imply that different models should be used, different information sets should be conditioned upon, and different assumptions are needed to justify conclusions.

## 1.4   A Guide to This Book

Our book spans fields within the social sciences, digital humanities, computer science, industry, and government. To convey our view on how to work with text as data in these disparate fields, we provide a different organization of our book. While most computational social science books organize the manuscript around algorithms, in this book we organize the book around tasks. We focus on tasks to emphasize what is different when social scientists approach text as data research. This also enables us to explain how the same algorithm can be used to accomplish different tasks and how validations for an algorithm might differ, depending on the goal at hand when applying that algorithm.

We organize our book around five key tasks: representation, discovery, measurement, prediction, and causal inference. Underlying this task-based focus is a set of principles of text analysis that we outline in Chapter 2. There, we explain our *radically agnostic* approach to text as data inference. We generally reject the view that there is an underlying structure that statistical models applied to text are recovering. Rather, we view statistical models as useful (and incomplete) summaries of the text documents. This view provides us with important insights into how to validate models, how to assess models that provide different organizations, and the role of humans within the research process.

In Part 2 we discuss selection and representation: the process of acquiring texts and then representing the content quantitatively. When selecting texts, basic principles of sample selection matter a great deal, even though there is a temptation to select content that is most conveniently available. When representing texts, we explain how different representations provide different useful insights into the texts and set the stage for future models in the book.

Part 3 introduces a series of models for discovery. By discovery we mean the use of models to uncover and refine conceptualizations, or organizations of the world. We show how a wide array of models can help suggest different organizations that can help researchers gain new insights into the world. We begin with methods used to uncover words that are indicative of differences between how two groups speak. These methods can be used to compare groups of documents—for example, legislators from two different political parties—or to help label categorizations inferred from other inductive methods. We then discuss some computer-assisted techniques for discovery, including models for partitioning data that exhaustively assign each observation to a single category. We then explain how clustering methods can be extended to admixture models, which represent each document as proportionally assigned to different categories. Finally, we describe methods for embedding documents into lower-dimensional spaces, which can shed light on underlying continuous variables in the data.

Part 4 describes our approach to measurement: assessing the prevalence of documents within a set of categories or assessing their location along a predetermined spectrum. We explain how to combine human judgment with machine learning methods to extend human annotations coded in a training set to a much larger dataset. When performing measurement, we explain how a discovery method can be repurposed to measure a category of interest. We include an extensive discussion of how to validate each of these measures, no matter what method produced them.

Building on the concepts and measures we have described, Part 5 explains how to apply the methods for prediction and causal inference. First, we describe how to use text as data methods to make predictions about how the world will be in the future. We discuss different types of predictive tasks and highlight how the threats to inference may vary with the setting. Next, we describe how to use the measures from texts as either the outcome or the intervention variable to make causal inferences. We explain the particular concerns that can emerge when text methods are used and provide a set of tools for assessing when a stringent set of assumptions is met.

## 1.5   Conclusion

There is immense promise with text as data research. With large amounts of data, complicated models, and custom measures, there is also the possibility of using these methods and getting the research wrong. Text is complicated and meaning is often

subtle. The risk is that if scholars overclaim on what text methods can do, they will undermine the case for using text methods.

Our book is intended as a guide for researchers about what is feasible with text as data methods and what is infeasible. We want to help readers learn about the immense set of tasks that text as data methods can help them accomplish. At the same time, we also want to help our readers to recognize the limits of text methods. We start out on this goal in the next chapter, where we articulate the basic principles that will guide our approach to text as data research.

# Index