

Big Data Applications in Historical Studies

Neil Eliason
Indiana University
Anderson, Indiana

ABSTRACT

As big data analytics progress in other fields, historians have begun to consider how they can apply these techniques to their studies. Various studies demonstrate potential benefits of big data approaches. However, care must be taken to keep big data results in the overall context of traditional scholarship and to utilize appropriate historical and technical expertise to avoid introducing inaccuracy and bias into findings.

KEYWORDS

i523, HID312, Big Data, History, Data Visualization, Inter-disciplinary

1 INTRODUCTION

1.1 Big Data

To date big data can claim numerous victories in a variety of fields, and promises more. Businesses such as Facebook and Netflix have built corporate empires off of the insights gathered from their big data, and physicists and biologists are learning what makes up the universe and ourselves via big data [1].

Despite all this, the concept itself is rather nebulously defined. A rough description is data with quantitative factors that require specialized techniques to utilize. The most commonly referenced big data factors are volume (amount of data), variety (number of data source types), and velocity (rate of data collection or input) known as “the three vs.” As these data factors become more extreme, to the point that traditional methods of data analysis fail, it becomes big data. While this definition is generally accepted, its application varies based upon the industry or field of study and often changes with developments in information technology [5].

The focus on big data arises partially from the phenomenon of data storage capabilities growing at a faster rate than data processing. This creates a situation where data can be economically stored, but not as economically processed, requiring specialized analytic techniques. As big data progresses through the storage, cleaning, analysis, and interpretation stages of the data life cycle, specialized approaches are required [1].

1.2 History of History

The historian’s labor has involved interacting with voluminous and varied data for centuries. Before computers, this process involved searching physical archives for relevant data, and manually copying and organizing it into useful information to be analyzed. Though this method can deliver deep insights, some data sets are too big to be studied in a manual fashion [7].

Around the mid-twentieth century, computers became sufficiently powerful and usable for historians to begin using them to process larger amounts of information. This facilitated a change towards a more quantitative approach to historical analysis and

a focus by some from tracing the rise and fall of political or ideological forces, to developing a more complete understanding of mundane topics, such as the family or economics.

As archives become digitized and accessible via the internet, the quantity of data available leads to an appeal to big data analytic methods [4]. The potential of unlocking significant connections and developing big picture historical insights at the scale of the growing digital archives of the world is alluring. This hope has driven the labor of many researchers towards developing more big data informed research methods and has directed funds of many institutions towards investments in data infrastructure. However, many are also concerned that the promises of big data are at best optimistic, and at worst hiding potential pitfalls to the historical process [7].

1.3 Thesis

Big Data Analytics have the potential to provide new insights to the field of historical studies. However, their application will differ due to the nature of historical data, and they will serve as an additional tool for the historian, rather than replacing more traditional approaches.

2 BIG DATA IN HISTORICAL STUDIES

2.1 Data Sources

It could be argued that history has had big data for some time, but that the lack of computational capability prevented it from being accessed on a large scale. As big data analytics mature, pressure develops to increase the data available for analysis by digitizing more archival material. This is evidenced not only by the familiar repositories of e-books, but also by archives of a variety of types, such as newspapers articles [7] or letters [4].

Sources for big data research consist not only of the content of documents in an archive, but also the bibliographical records. While originally designed to allow individual works to be located in an archive, historians have begun to study the bibliographical data themselves, an approach called distant reading. By looking at the data about a document, rather than the document’s content, societal or intellectual trends can be identified across large scale factors such as time or geography in a more comprehensive way. This approach has elicited some criticism that collections of bibliographical data are not complete enough to derive such large-scale conclusions. Still, considerable interest exists in targeting these data sets for historical analysis [10].

However, the data from these sources differs from that of other fields which utilize big data analytics. Historical data is not streaming the way that social media or smartphone sensors are. It is data which has already been collected, organized, and often times analyzed for a purpose defined by people from a different time and different needs/constraints from ourselves. This creates data

sets which are difficult to compare and often require considerable cleaning and reworking to be used in a larger framework. [4].

2.2 Analytics for Big Historical Data

Due to the natural reliance on documents in historical studies, text analytic techniques are the primary set of big data approach utilized by historians. Text analytics are a broad category of related algorithms and statistical techniques, such as artificial intelligence, machine learning, and natural language processing that attempt to extract specific information from the text and identify patterns and relationships within the body of data [7].

Artificial intelligence is “the ability of a digital computer or computer-controlled robot to perform tasks commonly associated with intelligent beings” [2]. In the context of historical research, this would include tasks such as extracting relevant content from sources or identifying relationships within the data. A specific type of artificial intelligence is machine learning, which consists of programs which change their actions autonomously in response to external input. Their ability to adapt allows them to do decision-making tasks, and thus can search through data sources in a more intelligent way to find relevant data [1]. Natural language processing is another artificial intelligence technique, which aims to create programs that can take human language, and make it machine readable [9]. Historians can use such programs to extract meaningful information from archival documents and prepare it for more further analysis and interpretation.

In order to interpret the results of big data analysis, visualization is critical. This is a challenge, as the large scale of the data makes striking a balance between a sufficiently big picture perspective without losing relevant details difficult. Many approaches attempt to utilize high resolution approaches to avoid losing important information [1]. This process is especially challenging in historical studies, as the data is often incomplete and may have inconsistencies which prevent assuming a uniform set of data. For this reason, historians often use visualizations to identify qualitative, rather than quantitative relationships in the data, to inform further inquiry [4].

2.3 Software Packages and Resources for Big Data History

A variety of software packages have been utilized to assist the process of translating raw data into historical insights, such as Tableau, Gephi, R, and ArcGIS. However, a limitation of these tools is their quantitative focus, which tends to exclude more qualitative approaches [4]. Some general qualitative analysis software has been applied to big data historical analysis, such as Google Fusion Tables and OpenHeatMap [10].

Some software has been developed to provide a more qualitative visualization tool set for researchers. For example Stanford University developed a software package called Palladio, designed to visualize connections in large scale historical data. Their approach focused on visualizations that encouraged exploring data, rather than creating statistical statements about it. Examples of this would be mapping connections between historical actors over geography or creating a visualization of the social network of a particular figure in history. They do not create statistical arguments, rather they

give a framework for understanding how the data are connected [6].

Another tool with a qualitative visualization focus is the Web Application for Historical Sentiment analysis on Public media or WAHSP. Its specific purpose is to conduct text analysis on the National Library of the Netherlands digitized newspaper collection, which contains around 100 million articles published in 1618 to 1995. It provides a number of useful analyses, such as word frequency cloud visualizations, detecting positive or negative sentiment related to certain terms, and Named Entity Recognition, which can identify people, places, events, etc. and then connect them into a relational or geographical framework. It also provides an interactive histogram where the resolution of the data can be adjusted to quickly move between a big picture and detailed data perspective. A derivative project is BILAND, which is a program developed by Utrecht University, that builds off of WAHSP’s analytical capabilities, but adapts them across the Dutch and German languages for comparative cultural studies [7].

Along with these data intensive tools specifically designed for historical studies, there are also resources to help the historian learn some of these methods. For example, The Programming Historian website provides a wide range of tutorials and lessons on how to use digital tools in historical studies. At the time of this writing there were 67 lessons available organized by their target stage of research, including lessons on using R, Python, Java, and GitHub for historical studies[8].

2.4 Insights from Big Historical Data

A number of studies have used these techniques to approach historical research from a big data perspective. Stanford’s Mapping of the Republic of Letters project sought to map the social network of Enlightenment thinkers who actively corresponded with each other. This was accomplished by utilizing big data analytics on the meta-data of these letters to see how these thinkers related temporally, geographically, and socially. Through the research process, the need for more qualitative approaches to visualization was recognized, and eventually led to the development of the Palladio tool set.

Their analysis revealed a number of interesting points. By mapping the social network of John Locke, they supported previous scholarly contentions that the Enlightenment culture was not homogeneously connected, but was made up of a number of subcultures which had thin social connections. Also, by analyzing Benjamin Franklin’s letters, they noted that despite his reputation as cross cultural traveler, the main hub of his correspondence was between the familiar British cultural hubs in Philadelphia and London [4].

Another study used the WAHSP tool to research attitudes found towards drugs in early 20th century newspapers. It found by using the word cloud analysis tool, that before 1924 drugs such as heroin and opium were discussed in the context of health, but after 1924 they were more associated with crime. Their analyses also noted that Dutch negative associations with opium influenced their perception of China and the Dutch East Indies Colonies.

The related tool BILAND was used by to study how the perceptions of eugenics differed in the Netherlands and Germany, requiring an application which could compare data across languages.

The aim was to study not only the direct conversations about this topic in both regions, but also to study implicit use of terminology which was influenced by the eugenics debate. Through word cloud analysis, the study found that in the mid 19th century, eugenics and concepts of genetic inheritance were used in a primarily medical or biological context. By the 1930s, the terms were utilized more in reference to race and law [7].

One study analyzed music bibliographical data from the British Library and the Repertoire International des Sources Musicales to explore how music was transmitted in Europe over time and geography. Their analytic methods were actually closer to traditional techniques, using a large amount of research assistants to perform repetitive tasks, and wrestling with the information in Excel spreadsheets, but used visualization approaches more congruent with big data. They had surprising results related to who were the prominently published composers during different time periods. For example, during the 1800s, relatively unknown composers are high in the frequency list, and famous composers such as Bach did not make the top 50 [10].

3 POTENTIAL ISSUES

While big data can provide some powerful and at times novel solutions to problems, there are also potential issues with its implementation. For example as digital algorithms make search and selection decisions, bias can be introduced into the research inadvertently by the program. This danger is aggravated by the level of transparency of the algorithm, and how well the researcher understands it. For example, when researchers utilize commercial search engines, such as Google scholar, the algorithms are not available, and thus the researcher does not know why data is being included or excluded. If recommender systems are utilized, the potential for bias increases, as the search engine is actively attempting to provide results which are based on its user profile. This could exclude opportunities for data which may challenge the researcher's perspective. The danger of biased analysis through ignorant execution of an automated search or analysis is present in any big data tool, such as those previously described [3].

In the context of historical studies, it is acknowledged that to use digital methods without expert knowledge of both the subject matter and the big data methodologies can lead to inaccurate conclusions [4]. However, this can be addressed using a number of strategies. For example, there are resources to help historians expand their technical abilities, giving them greater understanding and control over big data analytic methods [8]. Creating inter-disciplinary teams are also an effective way to address biased analysis. By allowing information technology and historical research experts to meet together to create research methods, they can avoid unintentional bias from misuse of algorithms and from a lack of knowledge of historical context. However, equally important is for the research team to keep a balanced perspective on the role of big data analytics applied in historical studies. These new methods cannot be done in a vacuum or be used to replace traditional human reading of the sources [4]. Though big data techniques have powerful possibilities, they cannot replace the role of the historian, who combines their historical knowledge and narrative creation, to provide context and meaning to the enormous bits of information from the past [7].

There are also a number of technical difficulties associated with using big data for historical analysis. The big data available to historical researchers has no guarantee of completeness or uniformity from which to make generalized claims. Large archives of records can only provide information about what people in the past chose to record or which records survived to our time. Thus, traditional methods are critical, and big data methods serve the purpose of confirming or challenging previous theories, or inspiring new veins of inquiry. History is an interpretative task, and big data analytics serve to better inform interpretation, not replace it. In addition, data often comes from different sources formatted for a variety of purposes. Thus for the historian, rather than dealing with large masses of unstructured data, the challenge is to reconfigure data which has already been organized, and often is at cross-purposes to a researchers objectives [4].

4 CONCLUSION

Big data analytics have attracted both interest and criticism from historians. Large digitized databases, effective text analytic techniques, and innovative qualitative visualizations provide fertile ground for a big data approach to historical analysis, which would allow for a more comprehensive analysis of large data sets, which would not be possible for the researcher. These techniques have already been applied to a variety of topics, yielding useful, if not incredibly surprising results.

As historians continue to explore new methods of big data research, it is important they do so from a position of historical and technical expertise, to prevent inaccurate and biased findings. The researchers' perspectives on big data analysis also needs to remain balanced, not ignoring the possibilities of the new techniques, but also not neglecting traditional research. Without traditional scholarship, big data has no external validation or historical context, thus making it's results inaccurate or meaningless.

ACKNOWLEDGMENTS

The author would like to thank Dr. Gregor von Laszewski for his support and suggestions to write this paper.

REFERENCES

- [1] C.L. Philip Chen and Chun-Yang Zhang. 2014. Data-intensive applications, challenges, techniques and technologies: A survey on Big Data. *Information Sciences* 275, Supplement C (2014), 314 – 347. <https://doi.org/10.1016/j.ins.2014.01.015>
- [2] B.J. Copeland. 2017. artificial intelligence (AI). Webpage. (01 2017). <https://www.britannica.com/technology/artificial-intelligence>
- [3] Malte C. Ebach, Michaelis S. Michael, Wendy S. Shaw, James Goff, Daniel J. Murphy, and Slade Matthews. 2016. Big data and the historical sciences: A critique. *Geoforum* 71, Supplement C (2016), 1 – 4. <https://doi.org/10.1016/j.geoforum.2016.02.020>
- [4] Dan Edelstein, Paula Findlen, Giovanna Ceserani, Caroline Winterer, and Nicole Coleman. 2017. Historical Research in a Digital Age: Reflections from the Mapping the Republic of Letters Project. *The American Historical Review* 122, 2 (2017), 400. <http://proxyiub.uits.iu.edu/login?url=https://search.ebscohost.com/login.aspx?direct=true&db=edsoaf&AN=eds0af.a29ec0ac934f1257030b477fa5986b1cff6def96&site=eds-live&scope=site>
- [5] Amir Gandomi and Murtaza Haider. 2015. Beyond the hype: Big data concepts, methods, and analytics. *International Journal of Information Management* 35, 2 (2015), 137 – 144. <https://doi.org/10.1016/j.ijinfomgt.2014.10.007>
- [6] Stanford Humanities and Design. 2017. Palladio. Visualize complex historical data with ease. Webpage. (2017). <http://hdlab.stanford.edu/palladio/about/>

- [7] Eijnatten Joris van, Pieters Toine, and Verheul Jaap. 2013. Big Data for Global History: The Transformative Promise of Digital Humanities. *BMGN: Low Countries Historical Review*, Vol 128, Iss 4, Pp 55-77 (2013) 128, 4 (2013), 55. <http://proxyiub.uits.uu.edu/login?url=https://search.ebscohost.com/login.aspx?direct=true&db=edsdoj&AN=edsdoj.6259f58bab47404485225cd4776fcf48&site=eds-live&scope=site>
- [8] Editorial Board of the Programming Historian. 2017. About the Programming Historian. Website. (10 2017). <https://programminghistorian.org/about>
- [9] Technopedia. 2017. Natural Language Processing (NLP). Webpage. (2017). <https://www.techopedia.com/definition/653/natural-language-processing-nlp>
- [10] Sandra1 Tuppen, Stephen2 Rose, and Loukia Drosopoulou. 2016. LIBRARY CATALOGUE RECORDS AS A RESEARCH RESOURCE: INTRODUCING 'A BIG DATA HISTORY OF MUSIC'. *Fontes Artis Musicae* 63, 2 (2016), 67 – 88. <http://proxyiub.uits.uu.edu/login?url=https://search.ebscohost.com/login.aspx?direct=true&db=llf&AN=114128249&site=eds-live&scope=site>

5 ISSUES

DONE:

Example of done item: Once you fix an item, change TODO to DONE

5.1 Assignment Submission Issues

DONE:

Do not make changes to your paper during grading, when your repository should be frozen.

5.2 Uncaught Bibliography Errors

DONE:

Missing bibliography file generated by JabRef

DONE:

Bibtex labels cannot have any spaces, _ or & in it

DONE:

Citations in text showing as [?]: this means either your report.bib is not up-to-date or there is a spelling error in the label of the item you want to cite, either in report.bib or in report.tex

5.3 Formatting

DONE:

Incorrect number of keywords or HID and i523 not included in the keywords

DONE:

Other formatting issues

5.4 Writing Errors

DONE:

Errors in title, e.g. capitalization

DONE:

Spelling errors

DONE:

Are you using *a* and *the* properly?

DONE:

Do not use phrases such as *shown in the Figure below*. Instead, use *as shown in Figure 3*, when referring to the 3rd figure

DONE:

Do not use the word *I* instead use *we* even if you are the sole author

DONE:

Do not use the phrase *In this paper/report we show* instead use *We show*. It is not important if this is a paper or a report and does not need to be mentioned

DONE:

If you want to say *and* do not use *&* but use the word *and*

DONE:

Use a space after . , :

DONE:

When using a section command, the section title is not written in all-caps as format does this for you

`\section{Introduction}` and NOT `\section{INTRODUCTION}`

5.5 Citation Issues and Plagiarism

DONE:

It is your responsibility to make sure no plagiarism occurs. The instructions and resources were given in the class

DONE:

Claims made without citations provided

DONE:

Need to paraphrase long quotations (whole sentences or longer)

DONE:

Need to quote directly cited material

5.6 Character Errors

DONE:

Erroneous use of quotation marks, i.e. use “quotes”, instead of ” ”

DONE:

To emphasize a word, use *emphasize* and not “quote”

DONE:

When using the characters & # % _ put a backslash before them so that they show up correctly

DONE:

Pasting and copying from the Web often results in non-ASCII characters to be used in your text, please remove them and replace accordingly. This is the case for quotes, dashes and all the other special characters.

DONE:

If you see a ffigure and not a figure in text you copied from a text that has the fi combined as a single character

5.7 Structural Issues

DONE:

Acknowledgement section missing

DONE:

Incorrect README file

DONE:

In case of a class and if you do a multi-author paper, you need to add an appendix describing who did what in the paper

DONE:

The paper has less than 2 pages of text, i.e. excluding images, tables and figures

DONE:

The paper has more than 6 pages of text, i.e. excluding images, tables and figures

DONE:

Do not artificially inflate your paper if you are below the page limit

5.8 Details about the Figures and Tables

DONE:

Capitalization errors in referring to captions, e.g. Figure 1, Table 2

DONE:

Do use *label* and *ref* to automatically create figure numbers

DONE:

Wrong placement of figure caption. They should be on the bottom of the figure

DONE:

Wrong placement of table caption. They should be on the top of the table

DONE:

Images submitted incorrectly. They should be in native format, e.g. .graffle, .pptx, .png, .jpg

DONE:

Do not submit eps images. Instead, convert them to PDF

DONE:

The image files must be in a single directory named "images"

DONE:

In case there is a powerpoint in the submission, the image must be exported as PDF

DONE:

Make the figures large enough so we can read the details. If needed make the figure over two columns

DONE:

Do not worry about the figure placement if they are at a different location than you think. Figures are allowed to float. For this class, you should place all figures at the end of the report.

DONE:

In case you copied a figure from another paper you need to ask for copyright permission. In case of a class paper, you must include a reference to the original in the caption

DONE:

Remove any figure that is not referred to explicitly in the text (As shown in Figure ..)

DONE:

Do not use `textwidth` as a parameter for `includegraphics`

DONE:

Figures should be reasonably sized and often you just need to add `columnwidth`

e.g.

```
/includegraphics[width=\columnwidth]{images/myimage.pdf}  
re
```