# Machine Learning and the Social Studies

Benjamin Charles Germain Lee, Ilene R. Berson, and Michael J. Berson

When we use digital assistants to process our search queries, generate driving directions, or answer our voice commands, advanced technologies gather available data and use this information to perform tasks that require decision making and problem solving. These examples of machine learning have become commonplace in our everyday lives. Machine learning refers to computer algorithms that demonstrate some form of intelligence akin to human cognition—for example, differentiating images of planes and trains, or identifying hate speech in social media posts. Much of our day-to-day technology is powered by machine learning. Given its astounding acceleration, most people do not realize that they are continuously encountering machine learning all around them.

Nonetheless, these technological advances can have problematic consequences that create ethical issues.[1] Recently, we witnessed the beneficial and dangerous outcomes associated with machine learning in social media contexts. Machine learning can forecast and disrupt the spread of disinformation, but it also has been used to generate entirely false news stories (i.e., deep fakes) that are indistinguishable from real news.[2] The tools used to counter disinformation also allow viral content to spread through social networks. When we come across unexpected but engaging posts in our social media feeds, machine learning algorithms have learned about our interests and preferences and curated content to capture our attention and focus. The side effects are technological echo chambers that propel malicious content into the public sphere and target users who are susceptible to extremist views.

We have a responsibility in the social studies to discuss machine learning, including its benefits and potential misuses as societies increasingly delegate complex decisions that affect the lives of citizens to artificial intelligence systems. We can best develop student competencies by using powerful machine learning tools as part of our instruction and pulling the curtain back on the technology that drives the innovation. This approach to media literacy in the social studies includes digital media learning experiences that delve into both the efficiencies and flaws to critically scrutinize technologies that have public policy implications.

## Newspaper Navigator

To integrate machine learning into the social studies classroom, students can consider how these innovations impact our ability to navigate digital collections, which are often limited not by the availability of digitized materials but rather by the ability to properly search them. Developments in machine learning have the capacity to unlock the ability to search over visual content, audio materials, webpages, and videos. As a case study, we discuss *Newspaper Navigator*, a project led by one of the authors in partnership with colleagues at the Library of Congress and the University of Washington to improve access to the visual content in 16+ million pages of digitized historic newspapers in *Chronicling America*.

*Newspaper Navigator* is a new free online tool for finding images and photos in the *Chronicling America* newspaper collection. The first goal of *Newspaper Navigator* was to use machine learning to extract photographs, illustrations, maps, comics, editorial cartoons, headlines, and advertisements from 16+ million pages of digitized historic newspapers. By extracting this visual content, along with captions, *Newspaper Navigator* enriches existing metadata for *Chronicling America*, allowing people to search more granularly than at the page level.

To train the machine learning model, *Newspaper Navigator* used thousands of crowdsourced annotations from the *Beyond Words* citizen history project,[3] which engaged the American public by asking volunteers to draw bounding boxes around the rich visual content in World War 1-era pages in *Chronicling America*.[4] By training a machine learning algorithm known as an object detection model to draw these bounding boxes, it became possible to automate this process across the entirety of *Chronicling America*. Figure 1 shows the visual content identified by this machine learning model on a sample page.

After running this process over all 16 million pages, the project team released the *Newspaper Navigator* dataset.[5] The dataset supports access via pre-packaged datasets of different visual content types

**Figure 1.** An example of visual content identified on a sample *Chronicling America* newspaper page by the *Newspaper Navigator* machine learning algorithm. The predicted category is shown in the top left of each predicted bounding box.
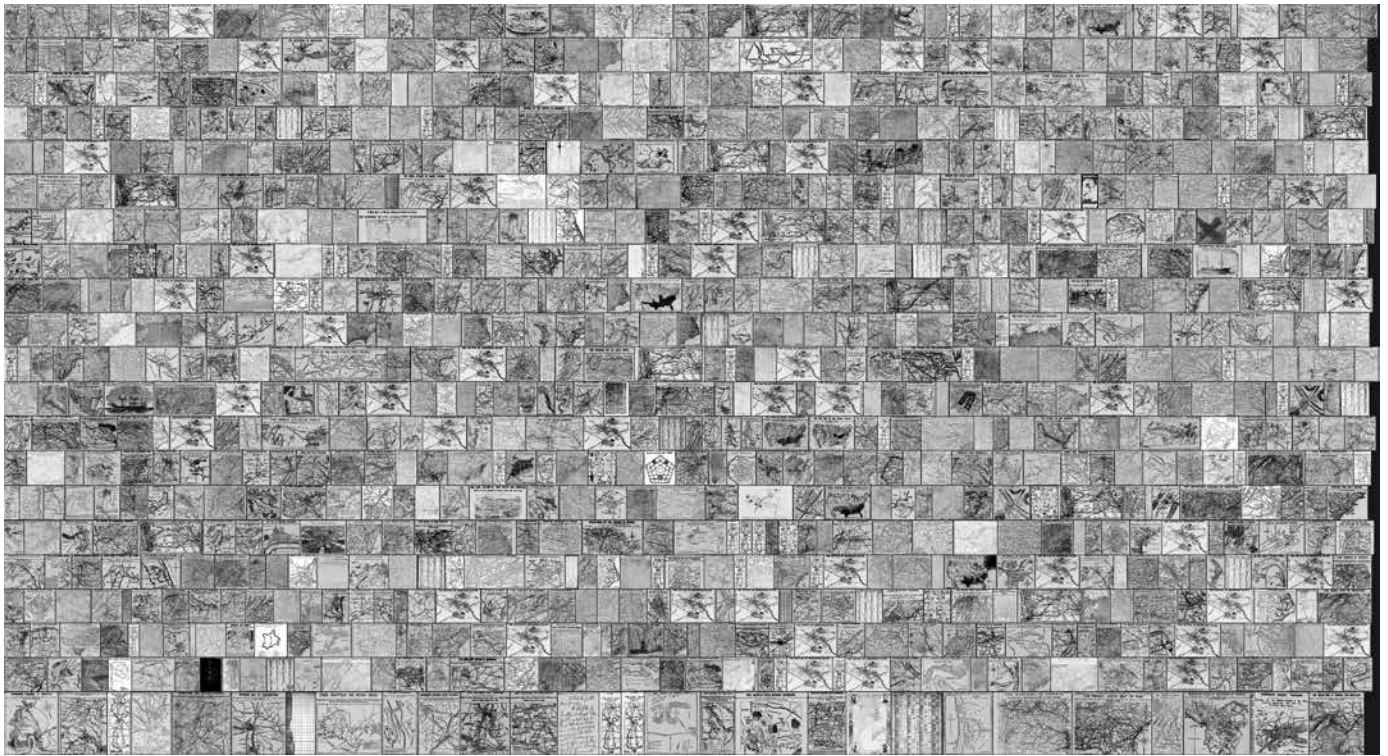
**Figure 2. A visualization of maps from the Civil War identified within the *Newspaper Navigator* dataset.**

## MACHINE LEARNING *from page 88*

separated by year, enabling the public to download and explore maps of the Civil War (as visualized in Figure 2) or editorial cartoons from World War 1. These pre-packaged datasets make it possible to explore many of the historic events and topics described in the *Chronicling America* topic pages at a larger scale.[6] These topic pages are already well-utilized among social studies teachers, and the *Newspaper Navigator* pre-packaged datasets can be used in tandem. The full dataset is in the public domain, making it suitable for all potential re-uses in the classroom.

With the *Newspaper Navigator* dataset constructed, the second goal of *Newspaper Navigator* was to reimagine how to search over the extracted photographs in the dataset. This step resulted in the *Newspaper Navigator* application for searching 1.5 million photographs published between 1900 and 1963.[7] In addition to supporting standard keyword search and filtering capabilities, the search application empowers visitors

to train a machine learning algorithm to search photos by visual similarity. Figure 3 illustrates how, based on user-selected examples of what the machine learning algorithm should and should not search for, the search application learns to retrieve desired types of photographs.

There are many applications for use in the classroom.[8] In social studies, this method of search offers students an opportunity to look at an event from multiple perspectives and unlock topics as far-ranging as emerging fashion trends or historic architectural styles. These user-defined content types can be as diverse as photographs of baseball players or sailboats or oval-shaped portraits. Students may trace the history of a topic over time or explore geographic differences in perspectives about an event. As they search for similar images, they may create a collection, examine the images for goodness of fit, and "retrain" the machine learning tool to improve its search performance. When the search process is complete, students may save the collection and have a URL generated to access it again or share with others.

These projects transform how students engage with historic sources and digital collections and provide a formative evaluation of students' historical thinking and visual literacy skills.

### Other Machine Learning Projects for Social Studies Teachers

Innovative applications of machine learning to digital collections are readily available for use in the social studies curriculum.[9] Some of our favorites are as follows:

- The Yale Digital Humanities Lab's PixPlot interface uses machine learning to explore 27,000 19th-century photographs.[10] Students may use the interface for visual analysis by examining clustered photos that share similar attributes (e.g., historical time period or content).

- Brian Foo's Citizen DJ project provides a freely accessible tool for exploring historic audio clips and moving images from the Library of Congress's collections.[11] Using
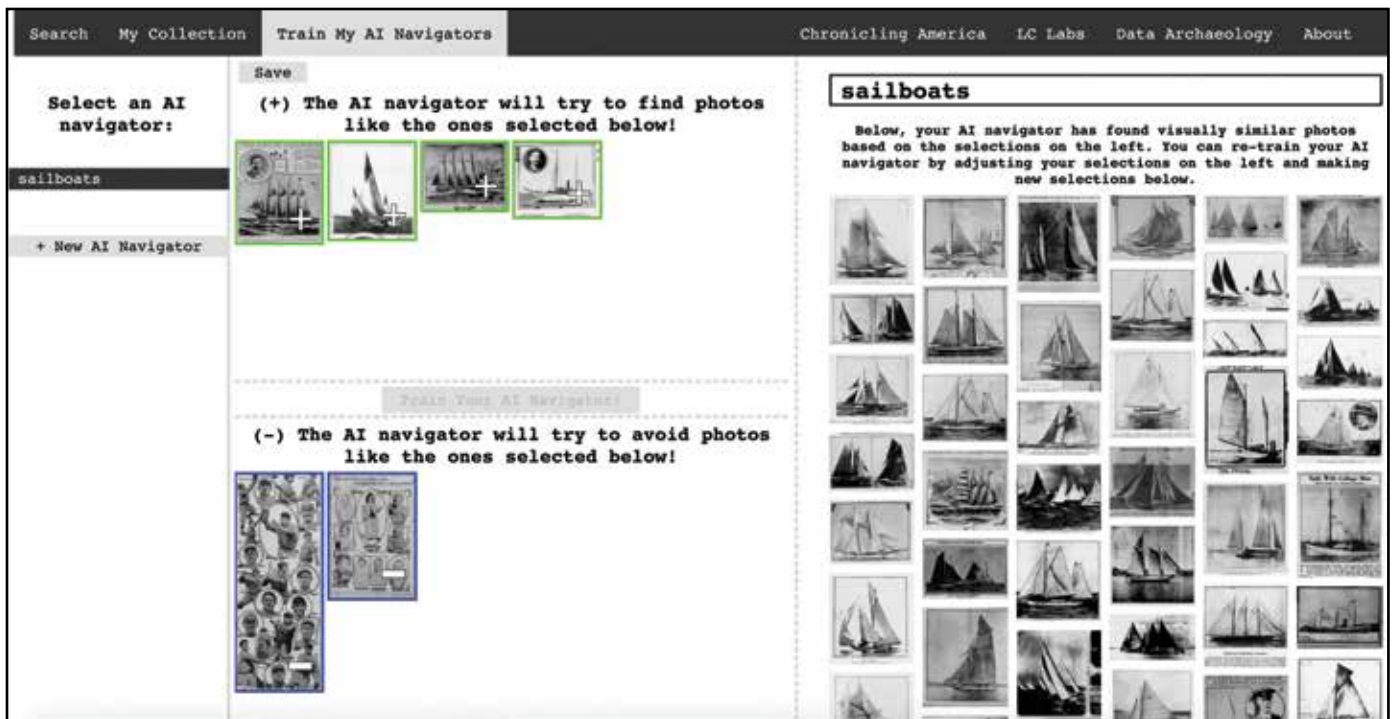
**Figure 3. A screenshot of the *Newspaper Navigator* search application showing how a visitor trains an AI navigator to retrieve photographs of sailboats.**

machine learning to generate samples, students may remix the recordings and create new music while exploring the historical context of the sounds. The materials include recordings of speeches, interviews, vaudeville acts, and other audio content.

- The "Dig That Lick" project navigates large databases of jazz solo recordings[12] and analyzes melodic patterns. Exploration of these musical patterns reveals how music is transmitted over time, and students can identify how music in the past has influenced other musicians across eras.

- The Civil War Photo Sleuth project uses artificial intelligence to identify photographs of unknown Civil War soldiers.[13] Teachers may use the interface to engage students in historical thinking skill development, using close inspection to corroborate matches.

- Google Arts & Culture Art Palette and MoMA's interactive archive of past exhibitions allow users to explore art with machine learning. Students may study color palettes in artwork and explore historical linkages across time.[14] They may also use the Google Art Selfie app to take a photo of themselves and have it compared with a vast archive of artwork to discover portraits that match their face and likeness.

## Machine Learning and Algorithmic Bias

Machine learning is already radically changing how we interact with digital collections. With these dramatic changes, we must also be aware of the ways in which machine learning can be harmful. Because most machine learning algorithms are trained on data labeled by people, these algorithms perpetuate marginalization by amplifying bias in the human-labeled data. This phenomenon is often called "algorithmic bias." Algorithmic bias can be found in a wide range of deployed machine learning

algorithms. For example, machine learning algorithms trained on text scraped from the Internet can learn to generate racist, sexist, and violent language. We must therefore be aware of how algorithmic bias affects our ability to search over cultural heritage collections when machine learning is used in the discovery process.[15]

To introduce students to algorithmic bias and explore how it affects the classification of information that may be used for decision making, students might experiment with Google's Teachable Machine Tool and try to build a classifier that distinguishes between cats and dogs when they are unknowingly provided with a biased dataset. When the classifier works better on cats than dogs, students have the opportunity to retrain their classifiers with their own new datasets.[16] Once they understand problems that can affect the accuracy of decision making with machine learning, students may discuss the potential social effects that occur when algorithms amplify biases and foster discrimination based on faulty assumptions by the human

programmers who design the technology. Secondary students may explore proposed legislation for algorithmic accountability or draft an algorithmic bill of rights to protect citizens.[17]

From algorithms that perform better on English than non-English languages to image recognition algorithms that perform poorly on images of people of color, the impacts of machine learning as applied to digital collections can potentially erase entire communities. Therefore, it is critical that students are educated about algorithmic bias and the harmful effects of machine learning whenever these tools are used in the classroom.

## Conclusion

In the coming years, machine learning will continue to shape how we navigate and access digital collections across different media, from documents to photos, audio, and webpages. In addition to improving metadata for collections, machine learning offers the ability to re-imagine how we perform searches and navigate materials, which present manifold opportunities in the classroom. With these exciting new methods of search comes the responsibility of understanding and teaching about the ways in which machine learning can also distort our understanding of our diverse cultural heritage. Social studies educators must not only promote an awareness of machine learning but must also address the broader media literacy issues that explicitly focus on the ways that machine learning can amplify inequality, racism, and injustice. 🌐

### Notes

1. See John Naughton, "From Viral Conspiracies to Exam Fiascos, Algorithms Come With Serious Side Effects," *The Guardian* (September 6, 2020), www.theguardian.com/technology/2020/sep/06/from-viral-conspiracies-to-exam-fiascos-algorithms-come-with-serious-side-effects

2. See John Villasenor, "How to Deal with AI-enabled Disinformation," Brookings, (November 23, 2020), www.brookings.edu/research/how-to-deal-with-ai-enabled-disinformation/

3. Learn more about *Beyond Words* here: https://labs.loc.gov/work/experiments/beyond-words/.

4. For more on crowdsourcing in social studies, see Ilene R. Berson and Michael J. Berson, "Crowdsourcing the Social Studies," *Social Education* 83, no. 2 (March/April 2019): 103–107.

5. You can access the dataset at https://news-navigator.labs.loc.gov. You can read more about the dataset construction at https://dl.acm.org/doi/10.1145/3340531.3412767.

6. Topic pages can be found at www.loc.gov/rr/news/topics/index.html.

7. The search application can be found at https://news-navigator.labs.loc.gov/search. An example of using the search function for the topic suffragist is discussed in Stephen Wesson, "Encouraging Student Exploration of Political Symbolism in Suffrage Cartoons," *Social Education* 84, no. 6 (November–December 2020): 384-386.

8. The podcast can be found at https://theprimarysourcepodcast.podbean.com/e/s1e3-the-newspaper-navigator-search-app-an-educators-view/.

9. A number of compelling examples of machine learning projects at cultural institutions are detailed in Ryan Cordell, "Machine Learning + Libraries," 2020, https://labs.loc.gov/static/labs/work/reports/Cordell-LOC-ML-report.pdf?loclr= blogsig.

10. Yale Digital Humanities Lab, "PixPlot: Visualizing Image Fields," https://s3-us-west-2.amazonaws.com/lab-apps/pix-plot/index.html?utm_source=dancohen&utm_medium=email.

11. "Designing for the (Citizen) | The Signal," webpage, https://blogs.loc.gov/thesignal/2020/09/designing-for-the-citizen-dj/. The interface can be found here: https://citizen-dj.labs.loc.gov/.

12. "Dig That Lick – Pattern Search | Search," https://dig-that-lick.hfm-weimar.de/pattern_search/

13. "Civil War Photo Sleuth," www.civilwarphotosleuth.com/; see also Vikram Mohanty et al., "Photo Sleuth: Combining Human Expertise and Face Recognition to Identify Historical Portraits," in *Proceedings of the 24th International Conference on Intelligent User Interfaces*, IUI '19 (New York, N.Y., USA: Association for Computing Machinery, 2019), 547–557, https://doi.org/10.1145/3301275.3302301.

14. Art Palette uses machine learning to find artworks based on a chosen color palette. See https://artsexperiments.withgoogle.com/artpalette/.

15. To learn more about these phenomena in the context of *Newspaper Navigator*, see Benjamin Lee, "Compounded Mediation: A Data Archaeology of the Newspaper Navigator Dataset" (September 1, 2020), https://hcommons.org/deposits/item/hc:32415/. For more on these effects in the context of libraries more generally, we recommend Thomas Padilla, "Responsible Operations: Data Science, Machine Learning, and AI in Libraries," OCLC, August 26, 2020, www.oclc.org/research/publications/2019/oclcresearch-responsible-operations-data-science-machine-learning-ai.html.

16. MIT Media Lab has designed a curriculum for middle school students to teach the ethics of artificial intelligence. See https://docs.google.com/document/d/1e9wx9oBg7CR0s5O7YnYHVmX7H7pnITfoDxNdrSGkp60/view

17. Sens. Cory Booker (D-NJ) and Ron Wyden (D-OR) proposed the Algorithmic Accountability Act of 2019 that would require companies to audit their algorithms for bias and discrimination. See www.wyden.senate.gov/imo/media/doc/Algorithmic%20Accountability%20Act%20of%202019%20Bill%20Text.pdf

**BENJAMIN CHARLES GERMAIN LEE** is a third-year Ph.D. student in the Paul G. Allen School for Computer Science and Engineering at the University of Washington. He was a 2020 Innovator in Residence at the Library of Congress, where he developed *Newspaper Navigator*. He can be contacted at bcgl@uw.edu.

**ILENE R. BERSON** is a Professor of early childhood education in the Department of Teaching and Learning at the University of South Florida. Her research explores early childhood social studies with a focus on the engagement of young children with digital technologies. She can be contacted at iberson@usf.edu.

**MICHAEL J. BERSON** is a Professor of social science education in the Department of Teaching and Learning at the University of South Florida and a Senior Fellow in The Florida Joint Center for Citizenship. His research focuses on technology in social studies education. He can be contacted at berson@usf.edu.