

DATA SCIENCE FOR INFORMED CITIZEN: LEARNING AT THE INTERSECTION OF DATA LITERACY, STATISTICS AND SOCIAL RESPONSIBILITY¹

Joachim Engel
Ludwigsburg University of Education, Germany, engel@ph-ludwigsburg.de

Focus Topics: AI and Data Science Education for Social Good, AI and Data Science Curricula and Implementation in School

Introduction

The information landscape is changing dramatically in the digital age due to the increasing availability of information via the internet, the widespread use of digital technologies, the abundance of data and easy access to data analysis tools. Digital media and the availability of data of sheer unlimited scope and magnitude change our access to information in radical ways. Emerging data sources provide new sorts of evidence, provoke new sorts of questions, make new sorts of answers possible and shape the ways in which evidence is used to influence decision making in private, professional and public life. In an increasingly data-driven world, social, societal and technological change requires new competencies. This expansion affects not only the professional world, but all of us. Innovation, social progress, and the well-being of our civil society require that people in science, business, politics, and society know how to evaluate and make sense of data to develop an informed understanding of our world and address pressing societal challenges with empirical insights and sound data-driven arguments.

At the same time, Big Data, with its possibilities for surveillance, manipulation, and control, poses serious problems for democracy and freedom (see, e.g., Helbing et al., 2017). The ability to assess the credibility of information and its sources has never been more important. The World Risk Report², published by the Swiss World Economic Forum Foundation in January 2024, sees misinformation and disinformation as the greatest threat to humanity over the next two years, ahead of extreme weather events, social polarization and armed conflict.

Algorithms drawing upon data are used to profile members of society and make crucial decisions which likely disproportionately impact those with less privilege and resources at their disposal. Amazon's model for sorting job applications³, for example, proved to be anti-women. Facebook's problems were first exposed by the Cambridge Analytica scandal, and the company continues to struggle with many ethical issues. Cathy O'Neil, in her book *Weapons of Math Destruction* (O'Neil, 2016), points out the dangers and injustices of using algorithmic models to determine credit scores, the price of insurance policies, whether someone should be paroled, or even what crimes police should investigate. Failure to learn how to understand, analyze and challenge data will result in citizens being in a continuously increasing position of informational disadvantage in relation to socio-political and commercial actors. Consequently, data literacy education needs to address a broad vision of data as social as well as technical assemblages. As consequence, data literacy and data science education cannot be reduced to learning technical mastery about algorithms, big data management and computing.

With all the promises of *Statistical Science to make a better world* (so a slogan of the International Statistical Institute), there are serious ethical concerns when more and more human activities are transcribed into data, quantified and analyzed (Van Es and Schäfer, 2017). Decisions taken by corporations and government agencies are increasingly data- and algorithm-driven, while the processes through which data are generated, communicated and represented are neither necessarily transparent nor devoid of negative effects (O'Neil, 2016). People are often unaware why, how or even that data about themselves are being collected, analyzed and 'shared' with additional parties (Dalton et

¹ Condensed version of Engel & Martignon (2024), <https://doi.org/10.37001/ripen.v14i3.3816>

² <https://www.weforum.org/publications/global-risks-report-2024/>

³ <https://www.ml.cmu.edu/news/news-archive/2016-2020/2018/october/amazon-scraps-secret-artificial-intelligence-recruiting-engine-that-showed-biases-against-women.html>

al., 2016). In an increasingly datafied society, data are often given the status of objective fact, despite its constructed, partial and biased nature.

Based on literature review, an analysis of needs to strengthen democratic values in the digital age and the reflection of own teaching practices, this talk aims to provide guidance on how to design data science education for informed citizens. After outlining the need for data literacy to be part of general education at any educational level and sketching some challenges for the emerging new field of data science education we present our concept of implementing some elements of data science and our objectives in classes for students preparing to be secondary school teachers. For this group no specific mathematical or technological background beyond high school is assumed, so the concept and goals may apply to any group of educated and informed citizens. We focus on how Data Science uses machine learning algorithms as one of its methodologies to analyze data, make predictions and automate decision-making processes.

Data Literacy – Challenges for the 21st Century Educators

Data science as a practical science has been conceived to address tangible problems in science, technology and society. Educating students in data science goes beyond teaching about algorithms, skills of manipulating data sets, selecting and applying appropriate analyses, and creating and interpreting visual representations of data. It also involves raising a critical understanding of how data are produced and how they can be used for particular purposes, including the role of context in interpreting data. It emphasizes developing an awareness for data ethics, and considering the implications for policy and society when powerful algorithms are used. Therefore, using real data is not enough, we need to teach data science to address real problems! Algorithms are not the goal of data science, they are an important tool. As a comparison: Physics is not about calculus but about understanding natural phenomena. Calculus is a tool for physics. Teaching data science for informed citizen and, e.g., future mathematics school teachers requires different contents than a course on data science for computer science, data engineers or statistics majors.

Data Science for Informed Citizens: Overarching Objectives

In this section, we discuss a personal account of what we believe every educated and informed 21st century citizen should know about machine learning and automated decision making. Our assertions are debatable, and you may find that some objectives are missing, may emphasize others more than we do, or disagree to a greater extent. However, the following contents have been implemented and evaluated in a series of classes for prospective secondary school teachers in mathematics and social studies. No specific mathematical or technological background beyond high school is assumed, so the concept and goals may well apply to any group of educated and informed citizens. Data Science is an interdisciplinary field that uses scientific methods, processes, algorithms, and systems to extract knowledge and insights from data with techniques from statistics, data analysis and machine learning (ML). Machine Learning is used in Data Science to make predictions or to classify data. It includes algorithms like decision trees, random forests, neural networks, and many others that automatically identify patterns and make decisions based on data.

Raise awareness where in daily life we encounter

Data Science products Data Science products and especially machine learning algorithms have become an integral part of our daily lives, often in ways many of us might not immediately recognize.

Teach awareness about data quality and data suitability

Data—the empirical basis for evidence-informed decisions and knowledge creation—are certainly preferable to anecdotes, wishful thinking, superstition, prejudice, or ideology. Yet data themselves are neither facts nor truth. Some authors consider data as models of reality. Data do not provide objective representations of the world. They might arise opportunistically, or as a result of conscious decisions someone made to research a particular topic. Data usually have been collected at costly expense, for a particular purpose and with a specifically chosen research design. They measure manifest variables in a particular way. They are the basis for constructing latent variables based on some

kind of model with a specific concept in mind. At a more complex level, one can ask why particular measures have been chosen, by whom, and for what purposes.

Collecting data is not a leisure activity but is laborious, sometimes tedious work that usually requires a lot of effort and financial resources. It serves someone's interest, and it is legitimate to question whose interest this is. Why have these data been collected? The data collected implicitly tell a story. Whose story is this? And whose story is this not? Critical or reflective questions about the methods used in surveys might include (but are not limited to):

- Are the measures (e.g., a questionnaire) well defined? Are the measures robust and appropriate for the purposes for which they are being used?
- Are metadata (i.e., detailed explanations of how variables were defined, sample characteristics etc.) available?
- Were the sampling procedures appropriate? Who is missing from the collected data?

Many studies in the social sciences are concerned with theories of causality; causality is associated with difficult philosophical challenges that go well beyond simple mantras such as “correlation does not imply causation.” However, when data come from observational studies, surveys, or archive data, and not from experimental studies, a reliable identification of cause and-effect relationships can be difficult to determine.

Beyond technical knowledge about processes of data generation, it is important that individuals are able to ask critical questions to assess the credibility and validity of any data, finding, or conclusion they encounter, both on technical and logical grounds. It is important to examine, from a critical perspective, narratives and interpretations of data, and the conclusions drawn from them

Teach awareness about biases of machine decisions

Machine learning algorithms can exhibit a lot of biases due to a variety of reasons, often reflecting issues in the data they're trained on, the design of the algorithm itself, or the broader societal and historical contexts in which they are developed and deployed. Here are several key factors contributing to bias in machine learning.

- Biased training data: Any bias in the training set will be amplified in the test set, leading to biased decisions. If the data used to train a machine learning model contains biases—either through underrepresentation or overrepresentation of certain groups, or through historical biases present in the data—the model will likely learn and perpetuate these biases.
- Algorithmic bias: The design of the algorithm itself can introduce bias. Some algorithms might be more prone to amplifying biases present in the training data. For example, algorithms that heavily penalize outliers might not perform well for minority groups that are underrepresented in the training data.
- Historical and societal context: The societal, historical, and cultural contexts in which data is generated often contain biases. Since machine learning models learn from past data, they can inadvertently learn and perpetuate these societal biases.

Addressing bias in machine learning is a great challenge that requires careful attention to the entire lifecycle of model development, from data collection and preparation through to deployment and monitoring. Strategies to mitigate bias include using more diverse and representative datasets, applying de-biasing techniques during model training, and continuously monitoring and updating models to ensure fair and equitable outcomes. The data that one uses needs to represent “what should be” and not “what is”. Otherwise, as in Amazon's hiring algorithm, the risk is high of underrepresenting and causes discrimination against a particular group of people. The validity of the algorithms needs to be evaluated when applied to various social groups. Furthermore, implementation needs some sort of mandated and enforced data governance to ensure a practice that is ethical with respect to the values of a free society (Shin, 2020).

Teach awareness about the impact of Data Science products on society

Creating critical awareness of the impact of machine learning on society is crucial to educating informed citizens who can contribute to the ethical development, deployment, and governance of these technologies.

In schools we can integrate discussions about the societal impacts of technology, including machine learning into the curriculum at various educational levels. This could range from simple lessons on digital literacy in elementary schools to more complex debates on ethics in high school and university courses. The ethics in machine learning encompasses legal, political social and economic dimensions. Therefore, an interdisciplinary approach that emphasizes cooperation with other fields such as ethics, philosophy, sociology or political science is appropriate. This approach ensures students not only learn how machine learning algorithms work but also understand their broader societal implications, including economic, social, and ethical dimensions.

An effective way is to engage students in projects where they have an opportunity to design, implement, or critique machine learning systems with ethical considerations in mind. This could include tasks such as creating a machine learning model taking into account potential biases, or developing guidelines for ethically deploying machine learning systems.

Teach some technological basics about machine learning

While algorithms are the tool and not the goal of Data Science, to appreciate the specific nature of machine learning, even students not majoring in computer science need to learn some of the technology about how a machine can do something we call “learning”. A good introduction are automated decision rules, represented by classification trees (Breiman et al., 1984). Trees are intuitive, simple to apply, easy to understand and give easily interpretable results. Decision trees created algorithmically from training data are simple and yet powerful tools capable of achieving high accuracy in many tasks while being highly interpretable. The “knowledge” learned by a decision tree appears as a hierarchical structure, a blueprint for decisions. This structure holds and displays the knowledge in such a way that even non-experts can immediately apply it.

Tree-based algorithms are an important method of machine learning which supports decision making, e.g., in medicine, finance, public policy and many more. Trees open doors to more advanced topics of Data Science and machine learning (e.g., Random Forests, Bagging and Boosting, as well as fundamental concepts such as training sets and overfitting).

However, instead of beginning with a computer algorithm that produces optimal trees, we suggest that students first construct their own trees, one node at a time, to explore how they work, and how well. This build it-yourself process is more transparent than using algorithms such as CART (Erickson and Engel, 2023). We believe it will help students not only understand the fundamentals of trees, but also better understand tree-building algorithms when they do encounter them. This may start completely unplugged. In a hands-on activity students receive data on cards. Each card has height, weight and various other measurements, including some irrelevancies (eye color) for a professional handball or football player. The students work to figure out how to predict the sport based on the other attributes. They came to see that they could cast their algorithm as a classification tree. Finally, they got previously-unseen cards to test their tree; this led, among other things, to discovering overfitting: a phenomenon where using the irrelevant variables resulted in an excellent, perhaps even perfect tree for the training data that was worse with new, test data.

Conclusions

Our main conclusions are:

Motivation problem definition and context

Data analysis must be motivated by a goal. And it must be embedded in a clear context in which it is to be applied and informed. A good Data Science application solves a well-defined problem or answers a specific question. This is the hard work that needs to be done before applying the automated tools. And it is one of the hardest things for students to learn and internalize at school and university.

Data provenance and metadata

The most sophisticated analysis is worthless if it is based on weak or questionable data. The context of the problem to be addressed is crucial for assessing the required relevance and quality of the data. Data analysis must not be conducted blindly, applied to data that is inappropriate or full of errors and gaps. Students must learn to document their data sources and their origin. And, more importantly, to be skeptical about the reliability of their data.

Human-machine interaction and decisions

Analytics must be a collaboration between human analysts and computer algorithms, with the algorithms serving as tools operated by humans. It is the human analyst who can adapt to changing circumstances, recognize the limitations of the model, understand the constraints of the data set, evaluate and correct, exclude or consider exceptional and deviant values, and understand the potential unintended consequences of a model that optimizes a criterion.

Ethics

Increasingly, ethical consequences of Data Science analysis are being uncovered. We must not rely on algorithms and must train our students to think and act ethically and apply these principles to their work. Students should learn to ask why an analysis is being performed and consider the ethical consequences of the answer. While many in the Data Science field view models as objective and unbiased, O'Neil (p. 21) defines models as "opinions embedded in mathematics." While the math gives the model the appearance of objectivity, in reality someone created the model and decided what data to use, what variables to include, what model form to use, and so on. A model is really an opinion that reflects both the bias of the modeler and the bias of the data itself. Those studying Data Science need to be sensitized to these ethical issues and trained in how to avoid bias and discrimination in models.

Problem solving

Sure, we also need to teach technical skills such as programming, machine learning algorithms and other big data topics. But that should not be the focus of a Data Science curriculum any more than calculus should be the focus of a physics curriculum. These are tools, and students should be good at them – but first they need to learn *why* and how to use them. The ultimate measure of success is solving the problem at hand by providing sustainable solutions that have tangible impacts.

Students should learn early in their education that Data Science is NOT about the tools! Data Science tools, no matter how powerful, are a “how”, not the “what”. Ultimately, it's not about knowing and using the tools well, it's about finding and using sustainable solutions to difficult problems. Otherwise, we shouldn't be surprised if the brightest minds we train use their brain power primarily to encourage other people to click on certain consumer ads rather than to use their knowledge to solve pressing social and societal problems.

References

- Atenas, J., Havemann, L., & Timmermann, C. (2020). Critical literacies for a datafied society: academic development and curriculum design in higher education. *Research in Learning Technology*, 28: 2468. <https://doi.org/10.25304/rlt.v28.2468>
- Breiman, L., Friedman, J.H., Olshen, R.A., and Stone, C.I. (1984). *Classification and regression trees*. Wadsworth.
- Dalton, C. M., Taylor, L., & Thatcher, J. (2016). Critical data studies: a dialog on data and space. *Big Data and Society*. 3 (1), 1–9. <https://doi.org/10.1177/2053951716648346>
- Engel, J. (2017). Statistical literacy for active citizenship: a call for data science education. *Statistics Education Research Journal*, 16(1), 44-49. <https://doi.org/10.52041/serj.v16i1.213>
- Engel, J., & Martignon, L. (2024). Data science for informed citizen: Learning at the intersection of data literacy, statistics and social justice. *Revista Internacional De Pesquisa Em Educação Matemática*, 14(3), 1-13. <https://doi.org/10.37001/ripen.v14i3.3816>
- Erickson, T., & Engel, J. (2023). What goes before the CART. Introducing classification trees with ARBOR and CODAP. *Teaching Statistics*, 45, S104–S113.
- Helbing, D., Frey, B., Gigerenzer, G., Hafen, E., Hagner, M., Hofstetter, Y., van den Hoven, J., Zicari, R. & Zwitter, A. (2017). Digitale Demokratie oder Datendiktatur. In: C. Könniker (Ed.), *Unsere digitale Zukunft*. https://doi.org/10.1007/978-3-662-53836-4_1
- O'Neil, C. (2016). *Weapons of Math Destruction: How Big Data Increases Inequality and threatens Democracy*. Crown Publishing Group.
- Richerich, A. (2018) *The Big Data Agenda: Data Ethics and Critical Data Studies*. University of Westminster Press, London. <https://doi.org/10.16997/book14>
- Van Es, K. & Schäfer, M. T. (Eds). (2017). *The Datafied Society: Studying Culture through Data*. Amsterdam University Press. <http://library.oapen.org/handle/20.500.12657/31843>