



BLOG @CACM

The *Communications* website, <https://cacm.acm.org>, features more than a dozen bloggers in the BLOG@CACM community. In each issue of *Communications*, we'll publish selected posts or excerpts.



Follow us on Twitter at <http://twitter.com/blogCACM>

DOI:10.1145/3575663

<https://cacm.acm.org/blogs/blog-cacm>

What Is Data Science?

Koby Mike and Orit Hazzan consider why multiple definitions are needed to pin down data science.



Why Is It Hard to Define Data Science?

December 2, 2022

<https://bit.ly/3B9tZxK>

If you ask a group of data scientists what data science is, you would probably hear different definitions. Indeed, although many attempts have been made to define data science, such a definition has not yet been reached. One reason for the difficulty to reach a single, consensus definition for data science is its multifaceted nature: it can be described as a science, as a research paradigm, as a research method, as a discipline, as a workflow, and as a profession. One single definition just cannot capture the diverse essence of data science. In this blog, we attempt to present the essence of each of these perspectives.

Data Science as a Science

Empirical science has been always about data. Kepler used data about the movement of the planets collected by Tycho Brahe to prove Copernicus' theory of the solar system. Was Kepler the first data scientist when he looked

for patterns and models in raw data? While Kepler used data to achieve insights, data science today is more than an empirical science. That is, *data science views the data itself as a natural resource and deals with methods for extracting value out of this data*.¹⁰ While science focuses both on understanding the world and on developing tools and methods to perform research, data science focuses on understanding data and developing tools and methods to perform research on *data*.¹¹

Data Science as a Research Paradigm

Data science also introduces a new scientific paradigm. The first scientific paradigm, established thousands of years ago, is *empirical science*, in which scientists describe natural phenomena. The second scientific paradigm, applied hundreds of years ago, is the *theoretical paradigm*, in which scientists build models of nature. The third scientific paradigm was introduced only several decades ago and is the *computational paradigm*, in which scientists simulate complex phenomena using algorithms and computers. The fourth scientific paradigm, according to Gray,⁷ is *data exploration*, in which data is captured or simulated, and then analyzed by scientists to infer new sci-

entific knowledge. Following Gray, the National Institute of Standards and Technology (NIST) claimed data science is the current evolution of the fourth paradigm and described data science as *"the conduct of data analysis as an empirical science, learning directly from data itself. This can take the form of collecting data followed by open-ended analysis without preconceived hypothesis (sometimes referred to as discovery or data exploration)."*¹³ In fact, this perspective views data science as the application of the grounded theory paradigm⁶ for quantitative research.

Data Science as a Research Method

Data science integrates research tools and methods taken from statistics and computer science that can be used to conduct research in various application domains, such as social science and digital humanities.

Drug discovery is one area that illustrates how machine learning is applied as a research method that includes *"target validation, identification of prognostic biomarkers and analysis of digital pathology data in clinical trials."*¹³

Another example is the application of machine learning methods in social science research. In such research, complex human-generated data, such as posts on social networks, are used to map social phenomena. Grimmer et al.⁸ reviewed current use of machine learning in social science research and stated, *"inclusion of machine learning in the social sciences requires us to rethink not only applications of machine learning methods but also best practices in the social scienc-*

es ...[machine learning] is used to discover new concepts, measure the prevalence of those concepts, assess causal effects, and make predictions. The abundance of data and resources facilitates the move away from a deductive social science to a more sequential, interactive, and ultimately inductive approach to inference.” As can be seen, data science is viewed in this case as a research method that transforms the research process from deductive to inductive, in line with the perspective on data science as a research paradigm presented here.

Data Science as a Discipline

Data science integrates knowledge and skills from several disciplines, namely computer science, mathematics, statistics, and an application domain. One way to present such a relationship is using a Venn diagram. Conway⁴ was the first to propose a Venn diagram for data science as a discipline; many other Venn diagrams were proposed for the discipline of data science following Conway.¹² Figure 1 shows our Venn diagram for data science.

Researchers recognize three levels of integration between two or more distinct disciplines: *multidisciplinarity*, *interdisciplinarity*, and *transdisciplinarity*.¹ **Multidisciplinarity** is the lowest level of integration. In multidisciplinary education, learners are expected to gain knowledge and understanding in each discipline separately. **Interdisciplinarity** represents a higher level of integration than multidisciplinary. In interdisciplinary education, after learners gain basic knowledge and understanding in each discipline separately, they are expected to understand the interconnections between the disciplines and to be able to solve problems that require applying different knowledge and methods from each discipline. In **transdisciplinarity**, boundaries among several disciplines transcend to create a holistic approach.

The gradual transition of data science from multidisciplinary to interdisciplinary introduces challenges concerning the definition of data science as a discipline that encompass the different bodies of knowledge, traditions, and cultures that originate in the different fields. Data science may eventually become a transdisciplinary domain.

Figure 1. The data science Venn diagram (the authors' version).

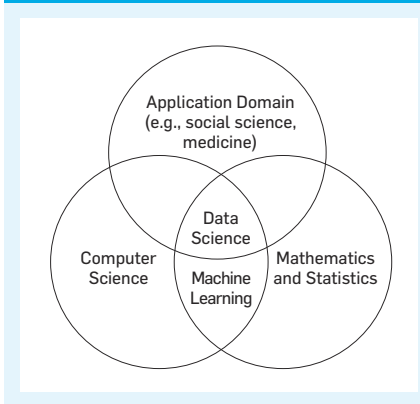
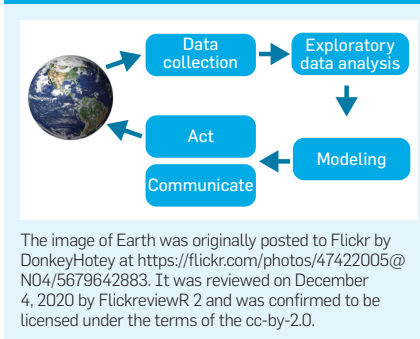


Figure 2. Data science workflow (the authors' version¹).



Data Science as a Workflow

The 2015 National Science Foundation (NSF) report, summarizing the NSF-sponsored workshop on data science education, introduced a definition of data science that reflects the perspective of data science as a workflow: “*Data science is a process, including all aspects of gathering, cleaning, organizing, analyzing, interpreting, and visualizing the facts represented by the raw data.*”² Data science is indeed commonly presented as an iterative workflow for generating value and data-driven actions from data (see Figure 2).

Data Science as a Profession

Irizarry⁹ proposed that the term *data science* was coined in order to improve communication between human resource recruiters in the industry and work applicants. According to Irizarry, “*As the demand in [sic] employees capable of completing data-driven projects increased, the term data scientist quickly became particularly prominent because it helped recruiters specify the type of employee they wanted.*”⁹

Accordingly, the definition of data science can be derived from the descrip-

tion of the profession of data scientist. And thus, for example, in their paper “*Data scientist: The sexiest job of the 21st century,*”¹⁵ Davenport and Patil describe the profession of data scientist as “*a high-ranking professional with the training and curiosity to make discoveries in the world of big data ... More than anything, what data scientists do is make discoveries while swimming in data. It’s their preferred method of navigating the world around them.*”

Conclusion

Data science is emerging fast and a single definition of it has not yet been reached. In this blog, we presented six facets of data science, each highlighting a different perspective of the field. This multifaceted nature of data science may partially explain the difficulties associated with the efforts to define it.

References

- Alvargonzález, D. Multidisciplinarity, interdisciplinarity, transdisciplinarity, and the sciences. *International Studies in the Philosophy of Science*, 25(4), 2011, 387–403. <https://doi.org/10.1080/02698595.2011.623366>
- Cassel, B. and Topi, H. *Strengthening data science education through collaboration: Workshop report 7-27-2016*. Arlington, VA, 2015.
- Chang, W. L., Grady, N., et al. *NIST big data interoperability framework: Volume 1, big data definitions*, 2015.
- Conway, D. The data science venn diagram. *Datist*, 2010. <http://www.dataists.com/2010/09/the-data-science-venn-diagram/>
- Davenport, T. H. and Patil, D. Data scientist: The sexiest job of the 21st century. *Harvard Business Review*, 90(5), 2010, 70–76.
- Glaser, B.G. and Strauss, A. L. *The Discovery of Grounded Theory: Strategies for Qualitative Research*. New York: Aldine de Gruyter, 1967.
- Gray, J. *EScience – A transformed scientific method*. http://research.microsoft.com/en-us/um/people/gray/talks/NRC-CSTB_eScience.ppt, 2007f
- Grimmer, J., Roberts, M. E., and Stewart, B. M. Machine learning for social science: An agnostic approach. *Annual Review of Political Science*, 2021 24, 395–419.
- Irizarry, R. A. *The role of academia in data science education*, 2020.
- Simberloff, D., Barish, B., Droegemeier, K., Etter, D., Fedoroff, N., Ford, K., Lanzetta, L., Leshner, A., Lubchenko, J., Rossmann, M., et al. *Long-lived digital data collections: Enabling research and education in the 21st century*. National Science Foundation, 2005.
- Skiena, S. S. *The data science design manual*. Springer, 2017.
- Taylor, D. Battle of the Data Science Venn Diagrams. *KDnuggets*. <https://www.kdnuggets.com/battle-of-the-data-science-venn-diagrams.html/>, 2016.
- Vamathevan, J., Clark, D., Czodrowski, P., Dunham, I., Ferran, E., Lee, G., Li, B., Madabhushi, A., Shah, P., Spitzer, M., et al. Applications of machine learning in drug discovery and development. *Nature Reviews Drug Discovery*, 18(6), 463–477, 2019.

Koby Mike is a Ph.D. graduate from the Technion's Department of Education in Science and Technology under the supervision of Professor Orit Hazzan. He is currently a post-doc at Bar-Ilan University. **Orit Hazzan** is a professor at the Technion's Department of Education in Science and Technology. Her research focuses on computer science, software engineering, and data science education. For additional details, see <https://orithazzan.net.technion.ac.il/>.

© 2023 ACM 0001-0782/23/2 \$15.00