

Embracing Data Science

Adam Loy

Department of Mathematics, Lawrence University

August 26, 2015

Abstract

Statistics is running the risk of appearing irrelevant to today's undergraduate students. Today's undergraduate students are familiar with data science projects and they judge statistics against what they have seen. Statistics, especially at the introductory level, should take inspiration from data science so that the discipline is not seen as somehow lesser than data science. This article provides a brief overview of data science, outlines ideas for how introductory courses could take inspiration from data science, and provides a reference to materials for developing stand alone data science courses.

1 Introduction

Statistics is running the risk of appearing irrelevant to many of today's undergraduate students. At first, this may sound absurd since the number of statistics degrees awarded at the undergraduate level approximately tripled from 2003 to 2013 [[Pierson 2014](#)]. I am not, however, referring to the popularity of the major, but rather the opportunities we as statistics educators have largely been missing. Rather than emphasizing that statistics is about “thinking with and about data” [[Cobb 2015](#), p. 3], we emphasize only part of the statistical thought process in our curricula. This issue is most pronounced at the introductory level. Consequently, we are failing to communicate what the field of statistics is about, making statistics seem irrelevant to many students, and perhaps even to many of our colleagues. In my experience many students who believe statistics is irrelevant are

interested in thinking with and about data, so where do they turn? Many are turning to data science.¹

2 What is Data Science?

In 2001, W. S. [Cleveland](#) outlined a plan for the discipline of statistics. This plan encouraged statistics departments to focus on the practice of data analysis, resulting in an altered field called data science. If one takes Cleveland’s view of data science, then it is a subset of statistics, essentially equivalent to applied statistics. Perhaps if every statistician, or at least a vast majority, had bought into Cleveland’s plan for the discipline this would be the case, but not all statistics departments/curricula resemble Cleveland’s ideal, so another definition has emerged.

Figure 1 shows an adaptation of Drew Conway’s venn diagram summarizing his view of data science [[Conway 2013](#)]. According to Conway, data science is the intersection of statistics, computer science², and domain knowledge. What is clear from figure 1 is that all three areas are necessary to define data science—e.g., without domain knowledge, we would just be talking about machine learning. Consequently, we cannot view data science as simply a subset of statistics, but rather it *utilizes a subset of statistics*. This intersection defining data science is what is often visible to undergraduate students in the media—for example, *The Upshot* section of *The New York Times* often presents data science products on a variety of topics. So it is important to understand data science in order to understand the background and expectations of our students.

Having defined what data science is, we must now turn to what a data scientist does. To better understand this, I have found the three steps of a data science project outlined by [Wickham \[2014a\]](#) to be a useful guide:

1. Collect data and questions.

¹I base this assessment on personal experience and observing the increasing popularity of data science training programs, such as the Coursera data science certificate run by Caffo, Leek, and Peng at Johns Hopkins and the Data Camp tutorials.

²Conway takes a narrower view, using the phrase “hacking skills,” but this view seems needlessly limiting.

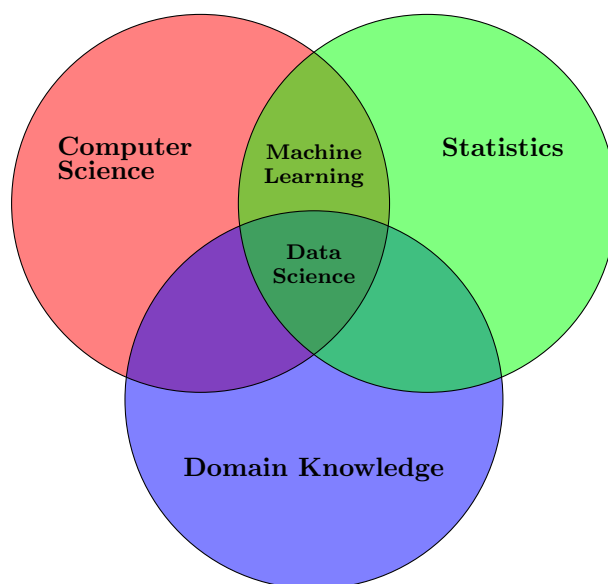


Figure 1: An adaptation of Drew Conway’s diagram explaining what knowledge and skills a data scientist needs. The original venn diagram was published on Drew Conway’s blog in 2010.

2. Analyze the data, using both exploratory and confirmatory methods.
3. Communicate the results in a way that convinces someone to “take action.”

At first, the above steps may sound identical to the steps in a applied statistics project, but I find statisticians and data scientists apply different weights to the tasks outlined above. For example, many statisticians downplay the importance of graphics in formal analyses believing that they are not rigorous enough, while data scientists fully embrace the power of graphical exploration. Further, many statisticians focus on formally reporting the results of their analyses, often in ways that are not accessible to a general audience. While the best applied statisticians are able to communicate their findings to a general audience, this seems to be a core requirement of a data science. Data scientists take communication to another level through the development of web apps and other data products, with the goal of appealing to a general audience.

While I have focused on the differences between the statistics and data science, there are far more similarities than differences as both disciplines are deeply rooted in the data analytic cycle. So why not take inspiration from data science when we design our statistics

courses?

3 How it can inform our teaching

Being mindful that today’s undergraduate students are familiar with data science projects and that they judge statistics against what they have seen will help us better teach statistics. Many students believe that statistics is irrelevant because we are not “rising to the level” of the data science projects they have seen. This is especially pronounced at the introductory level³ because we still mainly focus on the traditional pillars of statistics—data collection via surveys and experiments, and inference—rather than focusing on the entire data analytic cycle. Since both statisticians and data scientists think with and about data, we should take advantage of our students’ prior exposure to data science in our statistics curricula. In this section I outline a way to incorporate the elements of a data science project into an introductory statistics course. Such an introductory course should provide a trickle down benefit to the upper-level courses, assuming that students take an introductory course prior. I also provide a reference concerning the development of a stand alone data science course.

3.1 Introductory courses

The goals of an introductory course in any discipline are to (1) provide an overview of that discipline and (2) begin to develop a core knowledge base and skill set that is necessary for more advanced topics. I do not believe that many of the “traditional” introductory statistics courses achieve both goals, including those that I have taught. In fact, I think we often fail to achieve either goal because we try to cover too many modes of inference and some of us still emphasize by-hand calculations.

Rather than trying to cover so many inferential procedures we should take inspiration from data science and strive to provide an overview of the entire data analytic cycle. Focusing on the entire cycle will reveal the core knowledge base and skill set that should

³I consider the “traditional” introductory statistics course to be one that closely follows the AP curriculum [College Board \[2010\]](#). This includes the new randomization-based approaches to the course, e.g. [Lock et al. \[2013\]](#) and [Tintle et al. \[2015\]](#).

be emphasized. More specifically, by breaking down the steps of a data science project the necessary topics that are largely missing from the “traditional” introductory course emerge. These topics are core aspects of thinking with and about data, so we must incorporate them into our courses to stay relevant.

1. Collection

(a) Collect and refine questions

Too little time is spent questioning the data sets we use in class. While time constraints may be cited as a reason to discuss only “focused” analyses, a guided discussion of what questions could be answered using a specific data set will hone our students’ thought processes. This will allow them to question data sets outside of our classes, rather than only being able to rely on a data analytic “script.” Based on the recommendations from the American Statistical Association [ASA 2005; 2014], this change is already taking place, but we need to continue in this direction. One way to facilitate such discussions is to focus on case studies rather than textbook examples.

(b) Collect data relevant to the questions

Too often we provide students with the data directly, rather than having them collect or access data. In some courses, final projects provide students with their first opportunity to collect data. Rather than waiting until the end of the course, integrating data collection into the cycle will provide students with important tools for future analyses. In addition to traditional modes of data collection, it is essential to include some introduction to accessing web-based data.

2. Analysis

(a) Data manipulation

We must stop ignoring data manipulation. Data manipulation is often the most time consuming task in the data analysis cycle, so how can we justify it receiving little or no coverage in our introductory course? Discussing tidy data [Wickham 2014b] and the tools available to reshape data into the forms necessary for analyses are necessary additions to introductory courses.

(b) Visualization

Rather than going through a laundry list of exploratory graphics without motivation, we should incorporate graphics throughout the entire course. This will allow new graphics to be introduced on a just-in-time basis, and will avoid students thinking that graphics are unimportant or overly simplistic, as they seem to when we force the entire section into the first third of the traditional course. Further, we should refine the vocabulary we use in class by appealing to the grammar of graphics [Wilkinson 2006] and incorporate multivariate topics through the use of aesthetics and facetting. Kaplan [2015] provides an example of how this might be done at an appropriate level.

(c) Modeling

We already discuss statistical models in introductory courses, but why must we treat t -tests and simple linear regression models separately? A unified discussion of modeling from a linear model perspective, such as the one taken by Kaplan [2012], would help students focus on the thought process rather than the differences between formulas. Further, we need to discuss the inherent link between visualization and modeling in the data analytic cycle. This can easily be done if we model the thought process for our students during class and through assignments.

3. Communication

Many introductory statistics courses address communication in some way. Courses that require final projects seem to best develop a student's ability to communicate statistical results. However, we need to do more. Students are used to absorbing material online, so enabling them to communicate their findings online will have immediate impact. This can be done, for example, by teaching students how to build webpages using rmarkdown [Allaire et al. 2015] and knitr [Xie 2013], or basic web apps using Shiny [Chang et al. 2015]. These skills can be built incrementally at the end of each case study, and will provide very marketable skills, as well as the basic skills to perform reproducible research.

While adapting an introductory statistics course to include these topics requires a lot

of work, and numerous iterations to hone, the effort will be rewarded. After leaving a course where the entire data analytic cycle is examined I believe that students will be able to access new data sets and perform their own analyses in new situations. While less inferential material may be covered, students will understand how to manipulate and visualize data far better. (Many of the questions I field from former students relate to these topics, so they are certainly topics students encounter outside of our classes.) Additionally, less time will need to be devoted to developing basic data skills in later courses, adding room in the curriculum for more statistical topics. Finally, this approach provides a realistic overview of the field of applied statistics, and will help broadcast what statistics is about, even to those students we only see in an introductory course.

3.2 Teaching data science

In addition to changes in the introductory curriculum, the increasing visibility of data science opens the door for new courses. One possibility is to create a stand alone data science course. If there is not room in the curriculum for such a course, a hybrid course discussing more advanced statistical models and modern methods in a data science framework may be more feasible. I refer the interested reader to [Hardin et al. \[2014\]](#) for details from seven varieties of data science courses.

3.3 Conclusion

While statistics still seems to be increasing in popularity based on the number of degrees granted, it is running the risk of appearing irrelevant to a large population of students. If we take inspiration from data science and focus attention on the aspects of a data science project in our applied statistics courses, then we will show students that statistics is relevant while providing them with the tools to think with and about data, starting in the introductory course. This can only strengthen our applied statistics curricula and draw students into the discipline.

References

- Allaire, J., J. Cheng, Y. Xie, J. McPherson, W. Chang, J. Allen, H. Wickham, and R. Hyndman (2015). *rmarkdown: Dynamic Documents for R*. R package version 0.7.
- ASA (2005). *Guidelines for assessment and instruction in statistics education report*. Alexandria, VA: American Statistical Association.
- ASA (2014). *Curriculum Guidelines for Undergraduate Programs in Statistical Science*. Alexandria, VA: American Statistical Association.
- Chang, W., J. Cheng, J. Allaire, Y. Xie, and J. McPherson (2015). *shiny: Web Application Framework for R*. R package version 0.12.1.
- Cleveland, W. S. (2001). Data science: an action plan for expanding the technical areas of the field of statistics. *International statistical review* 69(1), 21–26.
- Cobb, G. W. (2015). Mere renovation is too little too late: We need to rethink our undergraduate curriculum from the ground up. *arXiv preprint arXiv:1507.05346*.
- College Board (2010). Statistics: Course description. Technical report.
- Conway, D. (2013, mar). The data science venn diagram.
- Hardin, J., R. Hoerl, N. J. Horton, and D. Nolan (2014, October). Data Science in Statistics Curricula: Preparing Students to "Think with Data". *ArXiv e-prints*.
- Kaplan, D. (2012). *Statistical modeling: A fresh approach* (2nd ed.).
- Kaplan, D. (2015). *Data Computing: An Introduction to Wrangling and Visualization with R* (Preview ed.). Project Mosaic Books.
- Lock, R. H., P. F. Lock, K. L. Morgan, E. F. Lock, and D. F. Lock (2013). *Statistics: Unlocking the power of data*. Wiley.
- Pierson, S. (2014). Bachelor's degrees in statistics surge another 20%. *AMSTAT News* (447), 27.

- Tintle, N., B. Chance, G. Cobb, A. Rossman, S. Roy, T. Swanson, and J. VanderStoep (2015). *Introduction to statistical investigations* (Preliminary ed.). Wiley.
- Wickham, H. (2014a). Data science: how is it different to statistics? *IMS Bulletin* 43.
- Wickham, H. (2014b). Tidy data. *Journal of Statistical Software* 59(10).
- Wilkinson, L. (2006). *The grammar of graphics*. Springer.
- Xie, Y. (2013). *Dynamic Documents with R and knitr*. Boca Raton, Florida: Chapman and Hall/CRC. ISBN 978-1482203530.

4 Acknowledgements

I would like to thank the participants of the Harnessing Big Data Workshop sponsored by the Associated Colleges of the Midwest. The conversations we had there made me think more specifically about my views on data science and its role in the undergraduate curriculum.