

## 1 ENSURING THAT MATHEMATICS IS RELEVANT IN A WORLD OF DATA SCIENCE

2  
3 Johanna S. Hardin (Pomona College) and Nicholas J. Horton (Amherst College)

4  
5 The recent growth of data science has been remarkable. Analysts now have rich data  
6 and powerful computational tools to help answer important questions. Examples of  
7 ways that insights can be wrangled from this information abound in diverse areas. This  
8 has led some to dub computational thinking (or fluency) as the "new literacy" on par with  
9 writing and quantitative skills. A major unanswered question relates to the role of  
10 mathematics in the training of future data scientists. How can we be sure that data  
11 science is on a firm mathematical and statistical foundation? In the article, we will  
12 consider what courses in mathematics would best prepare future data scientists.

### 13 14 **Background and brief history**

15  
16 Some institutions have responded to the development of data science by creating  
17 innovative new programs. At the University of California, Berkeley the Data 8  
18 introductory course (<http://data8.org/>) is now offered to a large proportion of incoming  
19 students, with connector courses on topics such as genomics, neuroscience, cultural  
20 data, social data, demography, smart cities, ethics, and social networks (as well as  
21 courses in statistics and mathematics). Many (most?) other four-year colleges and  
22 universities are responding with their own initiatives.

23  
24 While data science is often described as a new discipline, those in the mathematical  
25 sciences have been engaged with data science for decades. In a widely referenced call  
26 to action, Donoho (2017, in press), quotes noted statistician John Tukey from 1962 who  
27 presaged "an as-yet unrecognized science, whose subject of interest was learning from  
28 data, or 'data analysis' ". In his paper, Donoho describes the history of data science as a  
29 new field and speculates about a future that brings together statistics and machine  
30 learning by marrying computational and inferential methods. His proposed "Greater  
31 Data Science" (GDS) includes six main divisions (see Table 1).

32  
33 INSERT TABLE 1 AROUND HERE

34  
35 Table 1: David Donoho's Six Main Divisions for a "Greater Data Science" (Donoho,  
36 2017)

- 37 1. Data exploration & preparation: addresses the 80% (or more) of data wrangling  
38 needed prior to analysis
- 39 2. Data representation and transformation: including modern databases and special  
40 types of data
- 41 3. Computing with data: multiple environments, high-performance computing, and  
42 workflow
- 43 4. Data visualization and presentation: as a way to explore and present results in  
44 static or dynamic form
- 45 5. Data modeling: including both generative (stochastic model) and predictive  
46 (modern machine learning)
- 47 6. Science of data analysis: described as one of the most complicated of all  
48 sciences

### 49 50 **What mathematical preparation do future data scientists need?**

51  
52 What training is needed for data scientists to be able to extract meaning from data? This  
53 question is the topic for discussion by several working groups of the National Academy

54 of Sciences as well as a working group from the 2016 Park City Mathematics Institute.  
55 The potential for missteps, overgeneralization, and inferential errors abound. One of the  
56 challenges in training the next generation of students to 'think with data' is to ensure that  
57 they have sufficient background in the mathematical sciences to provide a firm  
58 foundation for their future work in data science.

59  
60 Unfortunately, many new data science programs have arisen that provide little or no  
61 formal preparation in the theoretical (mathematical, statistical and computational)  
62 underpinnings of this new field. While data science programs should appropriately focus  
63 on applications and practice, underlying many approaches is the use of modeling, a  
64 topic very familiar to the mathematical sciences, and abstraction, which underlies  
65 modern mathematics, statistics, and computational science. Practitioners need to  
66 understand when methods are applicable, where they are robust to underlying  
67 assumptions, and the potential for misbehavior. The danger is that students who skip  
68 out on math completely run the peril of "black box thinking", with no understanding of the  
69 *uncertainties* and *limitations* of models and algorithms. We argue that key concepts in  
70 statistics and mathematics undergird data science and that these essential aspects are  
71 needed as a foundation for data science. Additionally, we believe that mathematicians  
72 should take on the mantle of being directly involved in curricular decisions with respect  
73 to new data science programs.

74  
75 What kind of training in mathematics would be ideal for a future data scientist? It is not,  
76 we argue, the same training as would be ideal for a future mathematician. The proposal  
77 we outline below (two new courses on mathematics for data scientists) creates a path for  
78 integration of mathematics into data science. These new courses would not replace  
79 existing paths, since different preparation is needed for students who will be pursuing  
80 graduate degrees in mathematics.

81  
82 A gathering of computer scientists, statisticians, and mathematicians assembled at Park  
83 City Mathematics Institute (PCMI) during the summer of 2016 to propose guidelines for  
84 the discipline of Data Science (De Veaux et al, 2017). The group suggested that data  
85 science majors would indeed be well prepared by three semesters of calculus (including  
86 single and multivariable), Linear Algebra, Discrete Math, and Probability (in addition to  
87 several courses in statistics). They also noted, however, that such a course progression  
88 is *not feasible* for all students: it is *not realistic* for students to build a mathematical  
89 foundation that consists of such a long string of prerequisite courses *before* starting  
90 courses within their own "data science" curriculum (even if space could be made, the  
91 leakiness of lower-division pathways is a continuing problem, see  
92 <http://www.tpsemath.org/>).

93  
94 Project INGeniOuS (Investing in the Next Generation through Innovative and  
95 Outstanding Strategies, [http://www.maa.org/programs/faculty-and-](http://www.maa.org/programs/faculty-and-departments/ingenious)  
96 [departments/ingenious](http://www.maa.org/programs/faculty-and-departments/ingenious)) focused on ways that the mathematical sciences could help  
97 prepare the next generation of STEM students (at the same time that the mathematical  
98 sciences remained a vibrant choice for students). The joint report by the AMS, MAA  
99 (Mathematical Association of America), SIAM (Society for Industrial and Applied  
100 Mathematics), and the ASA (American Statistical Association) highlighted the  
101 importance of alternative curricular pathways and new approaches to teaching to ensure  
102 that the mathematical sciences are not left out of the growth of data science and other  
103 innovative interdisciplinary programs: "Curricula in the mathematical sciences  
104 traditionally aim toward upper-level majors' courses focused on theory. Shorter shrift is  
105 usually given to applications that reflect the complexity of problems typically faced in BIG

106 (Business/Industry/Government) environments, and to appropriate uses of standard BIG  
107 technology tools."  
108

109 How can the mathematics community respond to the challenge being posed by the  
110 growth of data science? We don't have all the answers, but we see the mathematical  
111 sciences as a key component of a vibrant and useful data science curriculum that  
112 provides students with a solid theoretical foundation. We suggest that the solution is to  
113 make changes to the mathematics and data science curricula to give future data  
114 scientists a glimpse into the power of mathematics and statistics for modeling and  
115 understanding a larger quantitative framework. Our fear is that the important  
116 mathematical foundational ideas will get lost if alternate pathways are not developed.  
117

### 118 **Mathematics preparation** 119

120 What then, is needed in terms of mathematical preparation? In order for students to be  
121 able to function effectively in the world of data science, we believe that that mathematics  
122 departments new to consider developing additional entry points as service courses.  
123

124 We propose two new courses - one discrete and one continuous (other approaches with  
125 similar pedagogic goals would also be natural to consider) that intertwine abstraction,  
126 modeling, and problem-solving. The idea of two new courses comes directly from the  
127 PCMI report:  
128

129         Mathematically speaking, the emphasis of an undergraduate data science  
130 degree should be on choosing, fitting, and using mathematical models. Because  
131 data-driven problems are often messy and imprecise, students should be able to  
132 impose mathematical [ideas] on [data science] problems by developing  
133 structured mathematical problem-solving skills. Students should have enough  
134 mathematics to understand the underlying structure of common models used in  
135 statistical and machine learning as well as the issues of optimization and  
136 convergence of the associated algorithms. Although the tools needed for these  
137 include calculus, linear algebra, probability theory, and discrete mathematics, we  
138 envision a substantial realignment of the topics within these courses and a  
139 corresponding reduction in the time students will spend to acquire them.  
140

141 Proposed New Course 1 (Mathematical Foundations I: Discrete Mathematics):  
142

143 The first proposed mathematics course formalizes the connections between  
144 mathematics and discrete model building (which leads naturally to more sophisticated  
145 topic and extensions in terms of continuous distributions, multivariate relationships, and  
146 causal inference). Combinatorial techniques can provide concrete pathways for  
147 explicitly conceptualizing models and their limitations. Linear algebra allows ideas of  
148 multivariate relationships, including independence. Many computer science  
149 departments teach a discrete course in their own departments. We suggest that those  
150 courses often focus more on algorithms as opposed to our suggestion that discrete  
151 models be used to conceptualize and model actual data and real world scenarios (and  
152 further develop the ability to problem solve using mathematics).  
153

154 The discrete topics suggested below would help the data scientist communicate about  
155 the multivariate problems they will inevitably encounter on a regular basis. Key discrete  
156 mathematical topics that would help a data scientist to model data effectively include:

- 157         • Linear algebra (ideas of independence / invertibility, Markov models and  
158             eigenvalues)

- 159 • Counting principles (understanding of first principles related to randomness)
- 160 • Computational (discrete) simulations associated with continuous models
- 161 • Graph theory (understand confounding, causal inference and analysis of network
- 162 data)

163

164 Proposed New Course 2 (Mathematical Foundations II: Continuous Mathematics):

165

166 A key aspect to modeling in data science is optimization. Part of what makes a model  
 167 appropriate has to do with its boundaries, maximal values, and sensitivity to parameter  
 168 choices -- all features that use mathematical optimization. In statistics, one foundational  
 169 method is to find parameter estimates by maximizing the relevant  
 170 likelihood. Alternatively, in other mathematical models, the goal might be nonlinear  
 171 state-space system identification. In both cases, a solid foundation of calculus,  
 172 differential equations, and numerical methods techniques will allow the data scientist to  
 173 solve the problem at hand. However, we argue that understanding how to find simple  
 174 minima and maxima (with ideas of local and global) acts as a vehicle for understanding  
 175 what optimization means at a fundamentally intuitive level. We recognize that the ideas  
 176 below are typically taught across many semesters. We are suggesting that much of the  
 177 content will be removed or taught differently so as to emphasize the critical mathematical  
 178 components necessary for data science. (For a model of such a course, see MATH 135,  
 179 Applied Calculus, taught at Macalester College to a large fraction of the undergraduate  
 180 population.)

181

182 To this end, the continuous mathematics course we suggest focuses on understanding  
 183 the continuous mathematical ideas necessary for problem solving. Some key topics to  
 184 be incorporated into such a course might include:

185

186

187

188

189

190

191

192

### 192 **The importance of computing**

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

Integrating computing into the mathematics curriculum not only serves students by  
 giving them computational skills, but additionally, technology in the mathematics

212 classroom allows students to understand the *mathematical theory* more completely. As  
213 the CUPM guidelines state:

214

215         In courses at all levels, substantial and realistic applications involve “messy”  
216 mathematics that makes calculation by hand onerous or infeasible. Using  
217 technology opens the door for students to set up solution strategies, justify their  
218 analyses, and interpret the results.

219

220 Using computational skills to simulate produces a deeper understanding of the model  
221 and complements analytic solutions. Additional computing will help develop better  
222 problem-solvers (and may yield additional mathematics majors drawn to the power and  
223 beauty of what they see in these courses).

224

225 While this article focuses on mathematical preparation, we believe that statistical  
226 preparation is also critically important. In recent years, the statistics community has  
227 taken on the challenge to improve their existing curriculum in order to ensure that  
228 statistical courses incorporate theoretical concepts, computation, and statistical practice  
229 (see for example the revised "Guidelines for Assessment and Instruction in Statistics  
230 Education [GAISE] College" report (ASA GAISE working group, 2016) (and the ASA  
231 revised "Guidelines for Undergraduate Programs in Statistics" (ASA Curriculum  
232 Guidelines working group, 2014). The latter report recommends that introductory and  
233 intermediate statistics courses (1) be an integral part of a data science curriculum, (2)  
234 incorporate reproducible research using statistical software (e.g., R Studio, Python  
235 notebooks, or GitHub), (3) use modern and relevant real data, possibly obtained through  
236 data scraping. While more work is needed by the statistics community, the article at  
237 hand primarily discusses the data science curriculum with respect to foundational  
238 courses in *mathematics*.

239

### 240 **Closing thoughts**

241

242 We see the world of data and modeling changing quickly. As mathematicians (and  
243 statisticians) we need to be proactive about what our disciplines have to offer.  
244 Mathematics will be better off if it is part of the solution. Data Science will be on a better  
245 foundational footing if it starts with mathematical first principles: abstraction and  
246 modeling. From students for many years, we understand at a visceral level how difficult  
247 it is for undergrads to grasp the benefits of generality and abstraction. Ensuring that  
248 they see the mathematical conceptual framework early and often will help make for  
249 better data scientists. In addition, abstraction is a key component of computer science  
250 and important linkages can be made.

251

252 We argue that mathematics needs to meet the growing data science community halfway  
253 so that the analysis and models leverage vital foundational mathematical concepts. If  
254 not, we run the risk that math will be left out. We have proposed one pathway to provide  
255 mathematical sophistication for beginning data scientists.

256

257 Our deliberately provocative suggestions, which build on the PCMI guidelines and the  
258 supplementary material therein, will not necessarily be easy to implement for many  
259 mathematics departments, given multiple competing interests and limited resources.  
260 However, we implore the community of mathematicians to take our suggestions  
261 seriously and engage in curricular discussions at their institutions so as to provide a  
262 strong theoretical framework to the world of data science and ensure that mathematics is  
263 not left behind. We look forward to working with our colleagues to develop multiple  
264 alternative approaches along the lines of those outlined by the Park City group in 2016.

265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291

## References

- ASA Curriculum Guidelines working group, Beth Chance, Steve Cohen, Scott Grimshaw, Johanna Hardin, Tim Hesterberg, Roger Hoerl, Nicholas Horton (chair), Chris Malone, Rebecca Nichols, and Deborah Nolan (2014), <http://www.amstat.org/asa/education/Curriculum-Guidelines-for-Undergraduate-Programs-in-Statistical-Science.aspx>
- ASA GAISE working group, Rob Carver, M. Everson (co-chair), John Gabrosek, Ginger H. Rowell, Nicholas J. Horton, Robin Lock, M. Mocko (co-chair), Allan Rossman, Paul Velleman, Jeffrey Witmer, and Beverly Wood (2016), "Guidelines for assessment and instruction in statistics education revised college report". [http://www.amstat.org/asa/files/pdfs/GAISE/GaiseCollege\\_Full.pdf](http://www.amstat.org/asa/files/pdfs/GAISE/GaiseCollege_Full.pdf).
- David Donoho (2017), "50 years of data science", *International Statistical Review*, <http://courses.csail.mit.edu/18.337/2015/docs/50YearsDataScience.pdf>
- Richard D. De Veaux, Mahesh Agarwal, Maia Averett, Benjamin S. Baumer, Andrew Bray, Thomas C. Bressoud, Lance Bryant, Lei Z. Cheng, Amanda Francis, Robert Gould, Albert Y. Kim, Matt Kretchmar, Qin Lu, Ann Moskol, Deborah Nolan, Roberto Pelayo, Sean Raleigh, Ricky J. Sethi, Mutiara Sondjaja, Neelesh Tiruvilumala, Paul X. Uhlig, Talitha M. Washington, Curtis L. Wesley, David White, and Ping Ye (2017), "Curriculum guidelines for undergraduate programs in data science", *Annual Review of Statistics and its Application*, DOI: 10.1146/annurev-statistics-060116-053930