

# A reproducible framework for undergraduate data science research

*Frederick J. Boehm and Bret M. Hanlon*

*8/20/2019*

## Contents

<b>1</b>	<b>Abstract</b>	<b>1</b>
<b>2</b>	<b>Introduction</b>	<b>1</b>
2.1	Our backgrounds . . . . .	2
<b>3</b>	<b>Methods</b>	<b>3</b>
3.1	Framework implementation . . . . .	3
3.2	Examples . . . . .	4
3.3	Relating to three ideas from Nolan and Temple Lang (2010) . . . . .	4
3.4	Broaden statistical computing to include emerging areas . . . . .	4
3.5	Deepen computational reasoning skills . . . . .	5
3.6	Combine computational topics with data analysis in the practice of statistics . . . . .	5
3.7	Develop and practice skills in reproducible research to promote open science . . . . .	6
<b>4</b>	<b>Outcomes</b>	<b>6</b>
4.1	Student Outcomes . . . . .	6
4.2	Scholarly outcomes . . . . .	6
<b>5</b>	<b>Discussion</b>	<b>6</b>
5.1	Benefits of our framework . . . . .	6
5.2	Critiques of our framework . . . . .	6
	References . . . . .	8

Last modified: 2019-12-05 16:58:16

## 1 Abstract

We design a mentoring framework to guide undergraduate researchers through individualized research projects in data science. Our framework involves research question formulation, data acquisition, data analysis and visualization, and presentation and communication of results. Our two honors students, whose projects serve as case studies for our framework, completed all components of the individualized research projects. We found that data science research skills, self-confidence in research ability, and professional interest in data science increased for both students. We describe our successes, lessons learned, and ideas for others to build similar frameworks.

## 2 Introduction

The need to analyze unprecedentedly large volumes of information combined with the development of faster and more powerful computers has fueled advances in data science methods for big data. Similar causes have led to a need for greater numbers of scientists with quantitative skills. In efforts to enhance training and

mentoring for trainees, we created a program that emphasizes many transferable skills that contribute to career success in data science.

We elected to work with social media data. This choice was deliberate. In making this decision, we recognized that social media data, such as tweets from Twitter, can be acquired with little cost and that there is growing research interest in social media in many social science disciplines, including political science, communication studies, and sociology. We also anticipated that our undergraduate trainees might be intrigued by the possibility of analyzing social media data, since many young adults use Facebook, Twitter, Instagram, and related sites.

Some social media data, including tweets from Twitter, are available through website application product interfaces (APIs). Twitter shares, via a streaming API, a sample of approximately one percent of all tweets during an API query time period (“Sampled Stream,” n.d.). Researchers have studied tweets for a variety of purposes, including inference of relationships and social networks among users (Lin et al. 2011); determination of authorship of specific tweets when multiple persons share a single account (Robinson, n.d.); and study of rhetoric in recruiting political supporters (Pelled et al. 2018; Wells et al. 2016). Recognizing the potential utility of tweets for data science research, we created a collection of tweets over time by repeated querying of the Twitter streaming API.

Nolan and Temple Lang (2010) argue for students to work with real data. Working with real data allows students to develop skill not only in statistical analysis, but also in data transfer from online sources, in data storage, and in using data from multiple file formats. In the case of Twitter data, tweets are stored in Javascript Object Notation (JSON) (“Consuming Streaming Data,” n.d.; “Introducing Json,” n.d.).

Mentoring in the work place and in higher education can have many benefits, including improving students’ development as thinkers and scholars, confidence in their own abilities, integration into the campus community, and interest in graduate training (Baker and Griffin 2010; Higgins and Kram 2001). A key component of our data science mentoring framework is the emphasis on using real data to answer real scientific questions. We believe that this process develops problem-solving skills that students will need in their future careers in data science. We encouraged the student to articulate a scientific research question, translate that question into quantitative and statistical terms, determine which data could be used to address the question, acquire the data, analyze data, visualize results, and communicate what they learned.

We provide guidance regarding selection of

## 2.1 Our backgrounds

During the time when we first implemented our framework, we served as early-career instructors in the statistics department at the University of Wisconsin-Madison. One of us (Hanlon) had prior experience in mentoring students, while the other (Boehm) had none. Our initial conceptualization of mentoring drew heavily on ideas we first encountered in professional development activities, including the Delta Program’s mentoring class (<https://delta.wisc.edu>) and Handelsman et al. (2005). Professor Erik Nordheim influenced our approach to and philosophy of teaching statistics. We studied with Professor Nordheim early in our teaching careers, and his emphasis on backward design and active learning continues to influence our teaching practices.

We both have experience in teaching undergraduate introductory statistics courses with enrollments over 100 students. Through our interactions with students in these classes, we’ve grown to value not only the ideas in a traditional introductory course, but also the need to prepare students with the essential skills needed for success in data science. Nolan and Temple Lang (2010) summarizes these skill sets in the following three goals:

1. broaden statistical computing to include emerging areas
2. deepen computational reasoning skills
3. combine computational topics with data analysis in the practice of statistics

To these three praiseworthy goals, we add a fourth:

84 4. develop skills in reproducible research to promote open science practices

85 We see the fourth goal as an equal with the first three from Nolan and Temple Lang (2010). Data scientists  
86 are uniquely positioned to promote open science practices, including the free sharing of data, code, and  
87 instructions for their use. The need for science to be more transparent and more reproducible elevate this  
88 goal to the level of the first three.

89 Below, we detail our methods for creating a reproducible framework for undergraduate data science research.  
90 We describe our results before concluding with lessons learned, things we could have done differently, and  
91 recommendations for future mentors who may use and extend our framework.

## 92 3 Methods

93 We designed and implemented a framework for mentored undergraduate data science research projects with  
94 big data. Below, we describe our initial framework and connect it to ideas from Nolan and Temple Lang  
95 (2010).

### 96 3.1 Framework implementation

#### 97 3.1.1 Research question formulation

98 Our mentored research framework begins with brainstorming scientific research ideas based on the student's  
99 interests. This enables us to craft a project that excites the student. With the results of brainstorming  
100 sessions, we (mentors and student together) formulate the most promising ideas into scientific hypotheses.

101 For the most appealing scientific hypotheses, we encourage the student to translate the scientific question  
102 into a statistical question that may be addressed with data. This is a crucial step in data science research  
103 question formulation. Skill in translating in both directions between scientific and statistical questions is a  
104 key communication skill that data science researchers offer.

#### 105 3.1.2 Data acquisition

106 We also incorporated data availability into our question formulation. We limited questions to those that  
107 could be studied with publicly available data. This practice also enabled reproducibility of our analyses, since  
108 students could share the URL from which they accessed data.

#### 109 3.1.3 Data analysis and visualization

110 After identifying research questions and publicly available data, the next step is to decide on informative  
111 data visualizations and quantitative analyses. Because both projects involved exploratory analyses of times  
112 series, we encouraged students to think about visualizations that might reveal relationships over time.

113 In the case of the event detection project, we e

#### 114 3.1.4 Presentation and communication of results

115 Students presented their research in a variety of settings. Each student presented at the annual undergraduate  
116 statistics poster session. We also encouraged them to present at the annual university-wide undergraduate  
117 research symposium.

118 In planning with students for poster and slide presentations, we (Hanlon and Boehm) emphasized the  
119 importance of succinctly stating the research question and its scientific context. After clarifying the

importance of the question, the student could proceed with explaining many of the elements that we've described above. Namely, the student would discuss the analyzed data and its acquisition while noting any shortcomings or biases of the data. For oral presentations, we suggested that students cautiously limit discussion of statistical methods, with the caveat that they prepare to answer detailed methodological inquiries during the question and answer session. Our students created powerful data visualizations for their projects. Their presentations also included their major results and future research directions.

In efforts to develop student written communication, we encouraged both students to prepare a written senior thesis document that detailed their research. In the senior thesis, we suggested that the students describe in rigorous detail their statistical methods. The rationale for this distinction, relative to the oral presentation, is that a reader doesn't have access to a question and answer session, while a poster session attendee may freely ask questions of the author.

## 3.2 Examples

Examples may help to demonstrate our approach to identifying a statistical research question. One of our students had interests in acquiring and using social media posts. We helped her in brainstorming ideas for research involving social media sources like Facebook and Twitter. Through this brainstorming, we recognized that she had a parallel interest in financial markets. Our student hypothesized that sentiment analysis of finance-related tweets might reflect trends in financial market index prices. On days when the market index prices increase, sentiment analysis of finance-related tweets might reveal more use of positive words, while days with decreasing prices might have more negative words in finance-related tweets.

A second student wanted to study tweets over time and entertainment events that garner lots of attention in social media. We encouraged this student to develop a strategy for event detection from tweets over time. The rationale is that a big entertainment event, such as the National Football League's Super Bowl game, might generate enough tweets that Super Bowl-related words would appear with high weights in results from latent Dirichlet allocation modeling of collections of tweets at distinct time points. We reasoned that Super Bowl-related topics might appear during the Super Bowl and vanish soon after the game's conclusion.

## 3.3 Relating to three ideas from Nolan and Temple Lang (2010)

We incorporate three key aspects that Nolan and Temple Lang (2010) identified:

1. broaden statistical computing to include emerging areas
2. deepen computational reasoning skills
3. combine computational topics with data analysis in the practice of statistics

Additionally, our projects gave students opportunities to develop and to practice skills in reproducible research. Given the growing imperative to document and share code to promote open science, we feel that this skill set equals in importance the three points above.

4. develop and practice skills in reproducible research to promote open science

Below, we describe how our framework enabled students to achieve competence in the four areas listed above.

## 3.4 Broaden statistical computing to include emerging areas

Our framework broadens statistical computing by including the emerging areas of social media data analysis, sentiment analysis, and topic modeling. Both students used Twitter tweets, which we accessed through a Twitter streaming API.

Our computational system for acquiring tweets involved several steps. We interacted with the API via the R package `twitterR` (Gentry 2015). We used the free Twitter streaming API that gave us access to approximately one percent stream of all tweets during the specified query time period. To ensure that we

collected tweets continuously, we used the linux tool `crontab` to execute our R script every five minutes. Each execution of the R script performed a single streaming API query for five minutes. Twitter’s streaming API, at the time of our data collection, enforced rate limits on the frequency and duration of queries. With the above settings, we continuously collected tweets.

### 3.5 Deepen computational reasoning skills

Our framework encourages students to deepen computational reasoning skills in several ways. First, they work with a variety of internet-based data to answer research questions. In the two example cases, our students collected tweets over time and gathered complementary data from other resources, including daily closing prices of stock market indexes. This gave students opportunities to think creatively about what data to acquire and how to use multiple data sources in a single cohesive project.

Second, the students worked with a variety of data structures. The Twitter streaming API returns tweets as JSON (Javascript Object Notation). Because distinct Twitter users may provide different pieces of profile information, there is variability in the structure of each tweet’s JSON. Additionally, tweet metadata fields may appear in any order (<https://developer.twitter.com/en/docs/tutorials/consuming-streaming-data>). Students needed to recognize this and to write code that accommodated these variations in tweet data structure. Additional variability in tweet structure arose due to changes in the API. The evolving nature of JSON tweet structure (<https://developer.twitter.com/en/docs/tweets/data-dictionary/guides/tweet-timeline>) required students to write flexible code that could incorporate newly introduced or deprecated metadata.

Students wrote R code to parse and organize tweet JSON []. They organized their R code into a package, and shared it on Github (<https://github.com/rturn/parseTweetFiles>). Each tweet’s JSON included required fields, and, possibly, some optional fields.

### 3.6 Combine computational topics with data analysis in the practice of statistics

Both mentored students combined computing with data analysis in the practice of statistics. They used a combination of latent dirichlet allocation topic modeling, sentiment analysis, and time series analysis to reach conclusions about real world data.

Both drew heavily on the collection of tweets. One student examined Standard and Poor’s 500 index daily closing prices over time. She also analyzed sentiments from each day’s stock market-related tweets to look for relationships between tweet sentiment and stock market prices.

Our other student focused on developing detection methods for social media events through topic modeling of tweets at different time periods. As a proof of principle, he fitted topic models to collections of tweets preceding, during, and following the National Football League’s Super Bowl game. He hypothesized that topics would evolve over time, with football-related tweets appearing during the football game and disappearing soon after conclusion of the game.

Both students analyzed tweets as texts. This first required them to write code to parse the JSON that the API returns. Once they had isolated the tweet text from its metadata, they parsed the tweet text into words for use in sentiment analysis and topic modeling. For the stock market project, they analyzed only those tweets that contained finance-related keywords. Sentiment analysis involved comparisons of tweet words to a dictionary that mapped words to sentiments. This yielded a net sentiment score for each tweet. They then treated tweet sentiment scores as a time series and compared them with daily stock market index closing prices.

The second student project involved latent Dirichlet allocation modeling of tweet words at distinct time points to detect social media events (Blei, Ng, and Jordan 2003). Latent Dirichlet allocation is a bayesian nonparametric method for modeling text corpora as the result of words chosen from topics.

### 3.7 Develop and practice skills in reproducible research to promote open science

With the goal of promoting transparency in our research, we encouraged students to use `git` for version control of their code and documents and to share their code via the website Github (<https://github.com>). One student also enrolled in Karl Broman's class on tools for reproducible research. This class features `git` and Github throughout its lectures and activities.

## 4 Outcomes

We applied the project framework to our mentoring of two students. Both engaged in 12 months of mentored research during their senior year of undergraduate studies in statistics.

### 4.1 Student Outcomes

We subjectively assessed student outcomes through conversations in our weekly student research meetings. Both students showed increases in confidence and ability to do data science research.

Both students secured positions in data science after graduation. One student later enrolled in a statistics graduate program, while the other pursued employment in health care analytics.

### 4.2 Scholarly outcomes

Our scholarly contributions include the `parseTweetFiles` R package on Github (<https://github.com/rturn/parseTweetFiles>) and presentations at conferences such as useR! 2016 and local poster sessions. Additionally, both students prepared end-of-project reports on their research.

We consider below

## 5 Discussion

### 5.1 Benefits of our framework

The student test cases for our framework demonstrated greater self-confidence and greater proficiency in data science skills over the course of the research projects. They used real-world data sources to address real scientific research questions. Additionally, they showed great interest in quantitative and data science careers. After graduation, one student immediately enrolled in statistics graduate training, while the other sought employment in health care analytics.

### 5.2 Critiques of our framework

From our current perspective, we offer a number of framework critiques and opportunities for improvement. Our measure of students' self-confidence in research ability was merely subjective. In future iterations of our framework, we would like to measure systematic and objective outcomes. One strategy for implementing this is to administer a survey, including questions from Vance et al. (2017), both before and after the mentored research project. We would use validated survey questions that focused on student beliefs about themselves, their skills, and their future careers.

One shortcoming of our initial framework was the relative lack of emphasis on best practices for computational reproducibility. This is one area that we would like to rectify in future mentoring activities. The University of

239 Wisconsin-Madison has periodically offered a semester course in best practices for computationally reproducible  
 240 research (<https://kbroman.org/Tools4RR/>). We especially see collaborative version control systems, such as  
 241 Git and Github, as essential tools for the modern data scientist.

- 242 1. assessment of data science skills
- 243 2. assessment of attitudes (pre and post survey??)

### 244 5.2.1 Framework development with backward design

245 In future research, we will continue to develop our framework for undergraduate data science research by  
 246 explicitly incorporating backward design principles (Wiggins and McTighe 2005). Following Wiggins and  
 247 McTighe (2005), we will identify desired results, determined acceptable evidence, and planned learning  
 248 experiences.

249 Before identifying desired results, we will prioritize topics from Nolan and Temple Lang (2010). Specifically,  
 250 we will assign all terms from Figure 1 of Nolan and Temple Lang (2010) into one of three categories:

- 251 1. worth being familiar with
- 252 2. important to know and do
- 253 3. enduring understanding

254 We've tabulated below the Nolan and Temple Lang (2010) terms for the current framework and its student  
 255 projects.

#### Prioritizing Key Terms from Figure 1 of @nolan2010computing

xxx	
Term	Circle
R packages	Enduring understanding
debugging	Enduring understanding
shell tools	Enduring understanding
reproducible computation	Enduring understanding
text editors	Enduring understanding
version control	Enduring understanding
file system concepts	Enduring understanding
text processing	Enduring understanding
regular expressions	Enduring understanding
EM	Important to know and do
MCMC	Important to know and do
Bayesian computation	Important to know and do
programming scope	Important to know and do
data structures	Important to know and do
portability	Important to know and do
authoring tools	Important to know and do
GUIs	Important to know and do
grammar of graphics	Important to know and do
composition	Important to know and do
linear algebra decompositions	Worth being familiar with
representation of numbers	Worth being familiar with
RNG	Worth being familiar with
optimization	Worth being familiar with
numerical algorithms	Worth being familiar with
efficiency	Worth being familiar with
parallel computing	Worth being familiar with
modeling language	Worth being familiar with
distributed computing	Worth being familiar with

compiled languages	Worth being familiar with
OOP	Worth being familiar with
symbolic math	Worth being familiar with
data bases	Worth being familiar with
I/O	Worth being familiar with
Flash	Worth being familiar with
HTTP	Worth being familiar with
XML	Worth being familiar with
SOAP	Worth being familiar with
SVG	Worth being familiar with
KML	Worth being familiar with
grid	Worth being familiar with
lattice	Worth being familiar with
event programming	Worth being familiar with
maps	Worth being familiar with
interactivity	Worth being familiar with
animation	Worth being familiar with
perception	Worth being familiar with
color	Worth being familiar with
raster/vector graphics	Worth being familiar with

---

Potential benefits of incorporating backward design ideas include clearer articulation of goals and better assessment of goal achievement.

We see our framework as one contribution to scholarship on improving data science training programs. Given the increasing economic need, in the USA and abroad, for data scientists and other researchers with quantitative training, we anticipate that our framework and its future iterations will continue to prepare students for data science careers by offering training in tangible and transferable analytic skills in the context of solving scientific questions.

## References

- Baker, Vicki L, and Kimberly A Griffin. 2010. "Beyond Mentoring and Advising: Toward Understanding the Role of Faculty 'Developers' in Student Success." *About Campus* 14 (6). Wiley Online Library: 2-8.
- Blei, David M, Andrew Y Ng, and Michael I Jordan. 2003. "Latent Dirichlet Allocation." *Journal of Machine Learning Research* 3 (Jan): 993-1022.
- "Consuming Streaming Data." n.d. <https://developer.twitter.com/en/docs/tutorials/consuming-streaming-data>.
- Gentry, Jeff. 2015. *Twitter: R Based Twitter Client*. <https://CRAN.R-project.org/package=twitter>.
- Handelsman, Jo, Christine Pfund, Sarah Miller Laufer, and Christine Maidl Pribbenow. 2005. *Entering Mentoring*.
- Higgins, Monica C, and Kathy E Kram. 2001. "Reconceptualizing Mentoring at Work: A Developmental Network Perspective." *Academy of Management Review* 26 (2). Academy of Management Briarcliff Manor, NY 10510: 264-88.
- "Introducing Json." n.d. <https://json.org>.
- Lin, Cindy Xide, Qiaozhu Mei, Jiawei Han, Yunliang Jiang, and Marina Danilevsky. 2011. "The Joint Inference of Topic Diffusion and Evolution in Social Communities." In *2011 Ieee 11th International Conference on Data Mining*, 378-87. IEEE.
- Nolan, Deborah, and Duncan Temple Lang. 2010. "Computing in the Statistics Curricula." *The American Statistician* 64 (2). Taylor & Francis: 97-107.



- 281 Pelled, Ayellet, Josephine Lukito, Fred Boehm, JungHwan Yang, and Dhavan Shah. 2018. “‘Little  
282 Marco,’ ‘Lyn’ Ted,’ ‘Crooked Hillary,’ and the ‘Biased’ Media: How Trump Used Twitter to Attack and  
283 Organize.” In *Digital Discussions*, 176–96. Routledge.
- 284 Robinson, David. n.d. “Text Analysis of Trump’s Tweets Confirms He Writes Only the (Angrier) Android  
285 Half.” <http://varianceexplained.org/r/trump-tweets/>.
- 286 “Sampled Stream.” n.d. <https://developer.twitter.com/en/docs/labs/sampled-stream/overview>.
- 287 Vance, Eric A, Erin Tanenbaum, Amarjot Kaur, Mark C Otto, and Richard Morris. 2017. “An Eight-Step  
288 Guide to Creating and Sustaining a Mentoring Program.” *The American Statistician* 71 (1). Taylor & Francis:  
289 23–29.
- 290 Wells, Chris, Dhavan V Shah, Jon C Pevehouse, JungHwan Yang, Ayellet Pelled, Frederick Boehm, Josephine  
291 Lukito, Shreenita Ghosh, and Jessica L Schmidt. 2016. “How Trump Drove Coverage to the Nomination:  
292 Hybrid Media Campaigning.” *Political Communication* 33 (4). Taylor & Francis: 669–76.
- 293 Wiggins, Grant, and Jay McTighe. 2005. *Understanding by Design*.