

## **Behavioral Data Science**

PSYC UN1930 (4 points)

*Instructor: Matthew R. Sisco*

[ms4403@columbia.edu](mailto:ms4403@columbia.edu)

*Office: Schermerhorn 403C*

*Office hours: XXX Xpm-Xpm; also available upon request*

Fall 2020

### **Syllabus**

**Bulletin Description:** This course covers the basic skills and knowledge needed to address psychological research questions using data science methods. Topics cover the full scope of a behavioral data science research project including data acquisition, data processing, and data analysis.

#### **Course Description**

This course covers the fundamental skills and knowledge needed for using data science methods to answer research questions in psychology. Topics will cover the full scope of skills and knowledge needed to complete a basic behavioral data science project including data acquisition (e.g. collecting data through APIs and web scraping), data processing (e.g. high performance computing and feature engineering), and data analysis (including machine learning, natural language processing, and advanced regression analyses). Discussion papers from the expanding literature at the intersection of behavioral science and data science will be examined and discussed. The papers will provide concrete examples of the techniques taught and will show the breadth of possible research designs. We will focus on papers addressing psychological research questions, and will also evaluate some relevant papers from other disciplines studying human behavior. The coursework will involve both empirical and methodological readings, and a series of skills-focused lab assignments (in the R programming language). As a core component of the class, students will design and implement their own behavioral data science research projects using methods taught in the course.

#### **Prerequisites**

- Science of psychology (PSYC UN1001) or similar
- AND Introduction to Statistics (PSYC UN1610 or equivalent) (you should be comfortable with simple linear regression models)

- AND Research methods (PSYC UN14xx or equivalent) or lab experience with research methods
- AND novice to intermediate R programming experience (you should be familiar with 80% or more of the basic R concepts and functions listed here: <https://rstudio.com/wp-content/uploads/2016/10/r-cheat-sheet-3.pdf>)
- AND Instructor permission

### **Course Objectives**

Through completion of this course, you will:

- a) Gain a broad theoretical and practical understanding of data science methods applicable to behavioral research
- b) Critically evaluate current approaches, research methods, and empirical papers in the field
- c) Constructively discuss relevant literature in class
- d) Develop your research communication skills, both oral and written
- e) Draw on course content to develop your own original research question and research proposal
- f) Design and implement a research project empirically evaluating your research question

### **Course Role in Departmental Curriculum**

This course is suitable for advanced undergraduates and postbacs majoring or concentrating in Psychology. Students pursuing a major in Neuroscience and Behavior will also receive registration priority. Students majoring in Economics, Sociology, or Political Science are eligible to enroll, however students majoring/concentrating in Psychology or Neuroscience and Behavior will be given registration priority.

Completion of this course satisfies the following psychology department requirement:

- Additional psychology elective course

Note: this course does not fulfill the psychology seminar requirement

### **Course Grading & Requirements**

- |     |   |
|-----|---|
| 10% | 1. Class participation (5% attendance, 5% class discussion participation) |
| 10% | 2. Thought papers   |
| 20% | 3. Weekly lab assignments   |
| 15% | 4. Project proposal (5% presentation, 10% written proposal)               |
| 35% | 5. Final project (10% presentation, 25% written paper)                    |

### **1. Class participation: 10%**

You are expected to attend and actively participate in every class. You should not only share your own thoughts on the discussion readings throughout the class, but also raise questions encouraging your peers to share theirs. Additionally, you will be expected to give your peers constructive feedback on their research proposals. Your participation will be evaluated after every class – as such, you will be penalized for any unexcused absences. Feel free to come see me anytime throughout the course to ask for feedback or suggestions regarding your class participation (or of course, to further discuss an idea that was raised in class). Participating in class can be more difficult for some students, and if that's the case, I encourage you to come see me at the beginning of the semester so that we can work out ways you can contribute. In these cases, later participation will be weighed more heavily to reward improvement.

### **2. Thought papers: 10%**

By 11PM the night before each class, you are required to submit a short thought paper on Canvas (roughly 150-250 words in length). The goal of these thought papers is to promote active reading and critical thinking, and to stimulate thoughts to discuss in class: you can raise theoretical or methodological questions related to the readings, share insights or comment on the implications of empirical findings, or relate the readings to previous class discussions. Try to integrate two or more readings into each thought paper. Be prepared to share your thoughts with your peers. These will not be formally graded but will be checked for completion/effort (each worth 1 point [those completed but with a clear lack of effort will receive half credit – note that greater length does not necessarily indicate greater effort]). Students can miss one thought paper during the semester at no penalty (but 1 extra credit point will be added to your final grade in the course if you complete all 11 of them).

### **3. Weekly lab assignments: 20%**

After each class meeting, students will have the remainder of the week to complete the related lab assignment. The lab assignments are due (submitted via Canvas) at the start of class the following week. There will be 10 graded assignments worth 20 points in total. The first assignment (“Assignment 0 – R programming basics review”) is optional but highly recommended for students with novice R programming experience. Students will receive 2 points for each on-time assignment, 1.5 points for each late assignment, and no points for incomplete assignments or assignments later than two weeks after the due date. The weekly lab assignments allow students to implement the techniques and methods discussed in class. There is some flexibility in the data and specifics of analyses students can use in the lab assignments. It is highly recommended that students use the lab assignments to prepare for the final project. It is perfectly acceptable to use analyses run in the lab assignments in the final project.

### **4. Project proposal: 15% (5% presentation, 10% written proposal)**

Mid-way through the semester, students will submit a 1-2 page project proposal and present a five minute summary of their proposed research project to the class. Each student is required to

schedule a meeting with me to discuss their ideas for a project proposal by the week before the day the proposals and presentations are due. At the end of the semester, students will summarize the analyses of their own data in a 10-15 minute class presentation and written paper based on the project proposed mid-way through the semester (more details on this final paper below).

Proposal structure:

- a. Introduction: Research question & brief literature review
- b. Proposed Data: Source of data and planned extraction procedures
- b. Proposed Method: Variables of interest and analysis procedure
- c. Predicted Results: Description of anticipated results (feel free to visualize)
- d. Discussion: Implications and limitations of predicted results
- e. References: 5+ references

### **5. Final project: 35% (10% presentation, 25% written paper)**

The research paper should (loosely) follow APA format with a brief introduction to the topic, a detailed methods section, a thorough results section, and a concise discussion. The paper should be between 1,500 and 3,000 words (not including references, figure captions, or tables). Final papers are due by submitting on canvas or paper copy on XXX at 5 pm and are worth 25 points.

Your papers will be graded based on thoroughness of the literature review (20%), integration of relevant and empirically valid methodology (20%), reasoning and implementation of the chosen analysis (25%), thoughtfulness of discussion (20%), overall presentation (grammar, spelling, APA formatting, etc.) (5%), and creativity and originality of the proposed idea (10%). The last grading criterion, creativity and originality, can be met by synthesizing ideas or approaches discussed in class to design your project, rather than simply replicating a discussed design with a minor variation.

In addition, you will present your research proposal (approximately 10-12 minutes including time for questions) to your classmates on the last day of class which is worth 10 points. I will discuss these presentations in more detail throughout the term.

### **Course Policies**

#### *Attendance:*

Absences will be excused with the presentation of proper documentation (i.e. doctor's or dean's note). Please inform me of the absence as soon as possible. You will still be responsible for completing the work due that particular class session. Each unexcused absence will result in 0.36 (5/14) points deducted from your class participation grade.

*Late work:* Unless excused by a Doctor's or Dean's note:

- Thought papers: Given that the purpose of thought papers is to prepare for the class discussion, you cannot submit a thought paper after class. Some leniency will be afforded for timing: half of your grade (0.5 points) will be deducted past the 11 PM deadline as long as it is submitted before the start of class.
- Lab assignments: 0.5 points off for late lab assignments.
- Project proposal: 1 point of your grade for the project proposal paper will be deducted per day the proposal is late. (Reminder the project proposal is worth 15 points total including the paper and presentation.)
- Final project: 1 point of your grade for the final project paper will be deducted per day the paper is late. (Reminder the final project is worth 35 points total including the paper and presentation.)

*Class Etiquette:*

Cell phones are not allowed to be taken out in class and should be kept on silent (not vibrate). Laptops or tablets may be used for anything course related. However, out of courtesy to your classmates and respect for your own learning, please refrain from using laptops or tablets for any other purpose.

*Students with Disabilities:*

If you are a student with a disability and have a DS-certified 'Accommodation Letter' please come to my office hours by the end of Week 2 to confirm your accommodation needs. If you believe that you might have a disability that requires accommodation, you should contact Disability Services at 212-854-2388 and [disability@columbia.edu](mailto:disability@columbia.edu).

*Academic Integrity:*

Columbia University Undergraduate Guide to Academic Integrity:  
<http://www.college.columbia.edu/academics/academicintegrity>

### **Faculty Statement on Academic Integrity:**

The intellectual venture in which we are all engaged requires of faculty and students alike the highest level of personal and academic integrity. As members of an academic community, each one of us bears the responsibility to participate in scholarly discourse and research in a manner characterized by intellectual honesty and scholarly integrity. Scholarship, by its very nature, is an iterative process, with ideas and insights building one upon the other. Collaborative scholarship requires the study of other scholars' work, the free discussion of such work, and the explicit acknowledgement of those ideas in any work that inform our own. This exchange of ideas relies upon a mutual trust that sources, opinions, facts, and insights will be properly noted and carefully credited. In practical terms, this means that, as students, you must be responsible for the full citations of others' ideas in all of your research papers and projects; you must be scrupulously

honest when taking your examinations; you must always submit your own work and not that of another student, scholar, or internet agent. Any breach of this intellectual responsibility is a breach of faith with the rest of our academic community. It undermines our shared intellectual culture, and it cannot be tolerated. Students failing to meet these responsibilities should anticipate being asked to leave Columbia.

### **Columbia College Honor Code:**

The Columbia College Student Council, on behalf of the whole student body, has resolved that maintaining academic integrity is the preserve of all members of our intellectual community – including and especially students. As a consequence, all Columbia College students make the following pledge:

We, the undergraduate students of Columbia University, hereby pledge to value the integrity of our ideas and the ideas of others by honestly presenting our work, respecting authorship, and striving not simply for answers but for understanding in the pursuit of our common scholastic goals. In this way, we seek to build an academic community governed by our collective efforts, diligence, and Code of Honor.

In addition, all Columbia College students are committed to the following honor code:

I affirm that I will not plagiarize, use unauthorized materials, or give or receive illegitimate help on assignments, papers, or examinations. I will also uphold equity and honesty in the evaluation of my work and the work of others. I do so to sustain a community built around this Code of Honor.

If found guilty of cheating or plagiarism, you will receive a zero for that assignment and be sent to the Dean ([www.college.columbia.edu/academics/disciplinaryprocess](http://www.college.columbia.edu/academics/disciplinaryprocess).) Please note that using code snippets from the internet IS acceptable, as long as you indicate where the code was copied from and provide a link to the public code source.

Citation should follow APA guidelines: <http://www.apastyle.org/>. If you have any doubt throughout the semester about how to cite something, or whether it would constitute as plagiarism, feel free to ask me.

### **Academic support services:**

Writing Center - <https://www.college.columbia.edu/core/uwp/writing-center>

Columbia Libraries - <http://library.columbia.edu/>

The schedule and materials listed below are subject to minor changes based on the progression of the class.

## Schedule and materials:

Week	Topics (Each class on the day a lab assignment is due will spend 30-45 minutes at the end of class reviewing the lab assignment)	Readings (To be read before the start of class each week, and discussed in class during the week they are listed for). (*Indicates reading is methodological and therefore will not be discussed in the same way as the research paper readings, rather it will be lectured on with questions taken.)	Assignments (due the following week)
1	<p>-Course overview -Ethics in data science research</p> <p>(if course is taught in a spring semester, week 1 will be spread over two weeks and include a statistics refresher)</p>	<p>Lazer, D., Pentland, A. S., Adamic, L., Aral, S., Barabasi, A. L., Brewer, D., ... &amp; Jebara, T. (2009). <a href="#">Life in the network: the coming age of computational social science</a>. <i>Science</i>, 323(5915), p721-723.</p> <p>Golder, S. A., &amp; Macy, M. W. (2014). <a href="#">Digital footprints: Opportunities and challenges for online social research</a>. <i>Annual Review of Sociology</i>, 40(1), p129-146.</p> <p>David Donoho (2017) <a href="#">50 Years of Data Science</a>, Journal of Computational and Graphical Statistics, 26:4, p745-766.</p> <p><i>Optional:</i> Lazer, D. &amp; Radford, J. (2017). <a href="#">Data ex Machina: Introduction to Big Data</a>. <i>Annual Review of Sociology</i>, p19-39.</p>	<p>Assn. 0: Basic review of R programming</p> <p>Thought paper</p> <p>(Reading page count: 42)</p>
DATA ACQUISITION			
2	<p>-Data acquisition overview -API data ingestion -Webscraping</p>	<p>Goel, S., Hofman, J. M., Lahaie, S., Pennock, D. M., &amp; Watts, D. J. (2010). <a href="#">Predicting consumer behavior with Web search</a>. <i>Proceedings of the National academy of sciences</i>, 107(41), p17486-17490.</p> <p>Ginsberg, J., Mohebbi, M. H., Patel, R. S., Brammer, L., Smolinski, M. S., &amp; Brilliant, L. (2009). <a href="#">Detecting influenza epidemics using search engine query data</a>. <i>Nature</i>, 457(7232), p1012-1014.</p> <p>Lazer, D., Kennedy, R., King, G., &amp; Vespignani, A. (2014). <a href="#">The parable of Google Flu: traps in big data analysis</a>. <i>Science</i>, 343(6176), p1203-1205.</p> <p>Ruths, D., &amp; Pfeffer, J. (2014). <a href="#">Social media for large studies of behavior</a>. <i>Science</i>, 346(6213), p1063-1064.</p> <p>Murphy, S. C. (2017). <a href="#">A hands-on guide to conducting psychological research on Twitter</a>. <i>Social Psychological and Personality Science</i>, 8(4), p396-412.</p>	<p>Assn. 1: API data ingestion</p> <p>Thought paper</p> <p>(Reading page count: 32)</p>
3	<p>-Open databases -Big surveys and experiments</p>	<p>Bond, R. M., Fariss, C. J., Jones, J. J., Kramer, A. D., Marlow, C., Settle, J. E., &amp; Fowler, J. H. (2012). <a href="#">A 61-million-person experiment in social influence and political mobilization</a>. <i>Nature</i>, 489(7415), 295-298.</p> <p>Bail, C. A., Argyle, L. P., Brown, T. W., Bumpus, J. P., Chen, H., Hunzaker, M. F., ... &amp; Volfovsky, A. (2018). <a href="#">Exposure to opposing views on social media can increase political polarization</a>. <i>Proceedings of the National Academy of Sciences</i>, 115(37), 9216-9221.</p> <p><a href="#">Resilient cooperators stabilize long-run cooperation in the finitely repeated Prisoner's Dilemma</a></p>	<p>Assn. 2: Webscraping</p> <p>Thought paper</p> <p>(Reading page count: 23)</p>

		<p>A Mao, L Dworkin, S Suri, DJ Watts - Nature communications, 2017, p 1-9.</p> <p>Kramer, A. D., Guillory, J. E., &amp; Hancock, J. T. (2014). <a href="#">Experimental evidence of massive-scale emotional contagion through social networks</a>. <i>Proceedings of the National Academy of Sciences</i>, 111(24), p8788-8790.</p>	
DATA PROCESSING			
4	<p>-Data processing overview</p> <p>-R in the cloud (linux, con jobs, FTP, SSH)</p> <p>-SQL, cluster computing</p>	<p>Strimas-Mackey, M. (2016). <a href="#">RStudio in the Cloud I: Amazon Web Services</a>. * (2 pages)</p> <p>Lane, R. (2019). <a href="#">Habanero – Getting Started</a> and <a href="#">R Job Examples</a>. HPC Cluster User Documentation. * (2 pages)</p> <p>Chapters 1 and 2 from Beaulieu, A. (2009). <a href="#">Learning SQL: Master SQL Fundamentals</a>. O'Reilly Media, Inc. p1-38*</p>	<p>Assn. 3: R in the cloud</p> <p>Thought paper (this week make it about ideas for your proposal since all papers are methodological)</p> <p>(Reading page count: 43)</p>
5	<p>-Statistical programming for long-run analyses (parallel programming, batch programming, best practices, GPU processing, scheduling)</p> <p>-Estimating additional variables (gender, ideology, age)</p>	<p>Wilson, G., Aruliah, D. A., Brown, C. T., Hong, N. P. C., Davis, M., Guy, R. T., ... &amp; Waugh, B. (2014). <a href="#">Best practices for scientific computing</a>. <i>PLoS biology</i>, 12(1), e1001745, p1-6.*</p> <p>Barberá, P. (2014). <a href="#">Birds of the same feather tweet together: Bayesian ideal point estimation using Twitter data</a>. <i>Political Analysis</i>, 23(1), 76-91.</p> <p>Matz, S. C., Schwartz, H. A., Menges, J. I &amp; Stillwell, D. J. (2019). <a href="#">Predicting Individual-level Income from Facebook Profiles</a>. PLOS ONE, p1-13.</p>	<p>Assn. 4: Long-run R programming</p> <p>Thought paper</p> <p>(Reading page count: 36)</p>
6	<p>-Monte Carlo simulations (for statistical evaluations and power analyses)</p> <p>-Bootstrapping</p>	<p>Muthén, L. K., &amp; Muthén, B. O. (2002). <a href="#">How to use a Monte Carlo study to decide on sample size and determine power</a>. <i>Structural equation modeling</i>, 9(4), 599-620.*</p> <p>Sigal, M. J., &amp; Chalmers, R. P. (2016). <a href="#">Play it again: Teaching statistics with Monte Carlo simulation</a>. <i>Journal of Statistics Education</i>, 24(3), 136-156.</p> <p><i>Optional:</i></p> <p>Ferron, J. M., Farmer, J. L., &amp; Owens, C. M. (2010). <a href="#">Estimating individual treatment effects from multiple-baseline data: A Monte Carlo study of multilevel-modeling approaches</a>. <i>Behavior Research Methods</i>, 42(4), 930-943.</p>	<p>Assn. 5: Monte Carlo simulations and boot strapping</p> <p>Thought paper</p> <p>(Reading page count: 43)</p>
DATA ANALYSIS			
7	<p>-Data analysis overview</p> <p>-Natural experiments (pre-post event analyses, IV regression, regression discontinuity)</p>	<p>Grimmer, J. (2015). <a href="#">We are all social scientists now: how big data, machine learning, and causal inference work together</a>. <i>PS: Political Science &amp; Politics</i>, 48(1), 80-83.</p> <p>Aral, S., &amp; Nicolaides, C. (2017). <a href="#">Exercise contagion in a global social network</a>. <i>Nature communications</i>, 8, 14753, p1-8.</p> <p>Sharma, A., Hofman, J. M., &amp; Watts, D. J. (2015, June). <a href="#">Estimating the causal impact of recommendation systems from observational data</a>. In <i>Proceedings of the Sixteenth ACM Conference on Economics and Computation</i> (pp. 453-470). ACM.</p> <p>Trochim, W. M. (2001). <a href="#">The regression-discontinuity design</a>. <i>International encyclopedia of the social and behavioral sciences</i>, 19, 12940-12945.*</p>	<p>No assn. or thought paper this week, instead prep for student proposals and presentations due next week</p> <p>(Reading page count: 37)</p>

8	<p>-Time series analyses (autocorrelation, stationarity, seasonality, unit roots, ARIMA)</p> <p>-Spatial analyses (spatial autocorrelation, spatial analysis programming, spatial regression modeling)</p> <p>-Student project proposals I</p>	<p>Chapters 1 and 2 of Pickup, M. (2014). <a href="#">Introduction to time series analysis</a> (Vol. 174). Sage Publications.* (49 pages)</p> <p>Chapter 1 of Ward, M. D., &amp; Gleditsch, K. S. (2008). <a href="#">Spatial regression models</a>. Sage.* (33 pages)</p> <p><i>Optional:</i> Pianta, S. &amp; Sisco, M.R. (2019). Is media coverage of climate change affected by changes in temperature? A six-year analysis of climate change media attention across 28 European countries. Working paper. (Paper will be uploaded to canvas) (10 pages)</p>	<p>Assn. 6: Time and spatial analyses</p> <p>Thought paper</p> <p>(Reading page count: 82)</p>
9	<p>-Public opinion estimation (with surveys, digital data, and multiple regression with post stratification (MRP))</p> <p>-Student project proposals II</p>	<p>Wang, W., Rothschild, D., Goel, S., &amp; Gelman, A. (2015). <a href="#">Forecasting elections with non-representative polls</a>. <i>International Journal of Forecasting</i>, 31(3), 980-991</p> <p>Beauchamp, N. (2017). <a href="#">Predicting and Interpolating State-Level Polls Using Twitter Textual Data</a>. <i>American Journal of Political Science</i>, 61(2), 490-503.</p> <p>Klašnja, M., Barberá, P., Beauchamp, N., Nagler, J., &amp; Tucker, J. (2017). <a href="#">Measuring public opinion with social media data</a>. In <i>The Oxford handbook of polling and survey methods</i>. (33 pages)</p>	<p>Assn. 7: Public opinion estimation</p> <p>Thought paper</p> <p>(Reading page count: 58)</p>
10	<p>-Machine learning (core concepts, popular classical and modern models, performance evaluation)</p> <p><i>*Guest presenter: Dr. Tal Golan (neural networks models)</i></p>	<p>Chapters 1, and 2.1-2.3 of James, G., Witten, D., Hastie, T., &amp; Tibshirani, R. (2013). <a href="#">An introduction to statistical learning</a>. New York: Springer, p1-37.*</p> <p>Gladstone, J. J.*, &amp; Matz, S. C.* (2019). <a href="#">Can Psychological Traits be Inferred from Spending? Evidence from Transaction Data</a>. <i>Psychological Science</i>, p1087-1096.</p> <p>Bhatia, S. (2019). <a href="#">Predicting risk perception: new insights from data science</a>. <i>Management Science</i>. (23 pages)</p> <p>Blumenstock, J., Cadamuro, G., &amp; On, R. (2015). <a href="#">Predicting poverty and wealth from mobile phone metadata</a>. <i>Science</i>, 350(6264), 1073-1076.</p> <p><i>Optional:</i> LeCun, Y., Bengio, Y., &amp; Hinton, G. (2015). <a href="#">Deep learning</a>. <i>Nature</i>, 521(7553), 436-444.</p>	<p>Assn. 8: Basic machine learning</p> <p>Thought paper</p> <p>(Reading page count: 75)</p>
11	<p>-Regression in machine learning (automated model selection, LASSO regression)</p>	<p>Tibshirani, R. (1996). <a href="#">Regression shrinkage and selection via the lasso</a>. <i>Journal of the Royal Statistical Society: Series B (Methodological)</i>, 58(1), 267-288.*</p> <p>Reece, A. G., Reagan, A. J., Lix, K. L., Dodds, P. S., Danforth, C. M., &amp; Langer, E. J. (2017). <a href="#">Forecasting the onset and course of mental illness with Twitter data</a>. <i>Scientific reports</i>, 7(1), 13006. (9 pages)</p> <p>Bedi, G., Carrillo, F., Cecchi, G. A., Slezak, D. F., Sigman, M., Mota, N. B., ... &amp; Corcoran, C. M. (2015). <a href="#">Automated analysis of free speech predicts psychosis onset in high-risk youths</a>. <i>npj Schizophrenia</i>, 1, 15030. (6 pages)</p> <p>Yarkoni, T., &amp; Westfall, J. (2017). <a href="#">Choosing prediction over explanation in psychology: Lessons from machine learning</a>. <i>Perspectives on Psychological Science</i>, 12(6), 1100-1122.</p>	<p>Assn. 9: ML with regression</p> <p>Thought paper</p> <p>(Reading page count: 59)</p>
12	<p>-Natural Language Processing I (word counting, feature extraction, NLP and machine learning, ensemble models)</p>	<p>Tausczik, Y. R., &amp; Pennebaker, J. W. (2010). <a href="#">The psychological meaning of words: LIWC and computerized text analysis methods</a>. <i>Journal of language and social</i></p>	<p>Assn. 10: Basic natural language processing</p>