

Psych/CSCI 626: Text as Data

Morteza Dehghani

Fall 2021

E-mail: mdehghan@usc.edu
Office Hours: Mon 10-12
Office: SGM 607

Web: cssl.usc.edu
Class Hours: Th 2-6pm
Class Room: BCI 266

Course Description

Text as Data focuses on applications of natural language processing, guided by psychological theories, for identifying various social and cognitive properties evident in human related big data. In this course, we will survey state-of-the-art techniques, and applications of such techniques, for investigating various aspects of human cognition. The intended audience for this course is psychology and computer science PhD students, and more broadly graduate students in social sciences, who are interested in using machine-learning techniques for analysis of data. Also, this course may be of interest to PhD students in communications and the business school.

Learning Objectives

This course is designed to survey current state of research in automated analysis of language within the domain of psychology. It should be noted that the purpose of this class is not to teach text analysis, nor social psychology, but to survey how the established methods are used within the social sciences. Optional reading material will be provided for students unfamiliar with topics discussed.

- **Prerequisite(s):** Instructor permission
- **Recommended Preparation:** For Non-Engineering majors: Psych 625 or a similar course, for Engineering students: CSCI 544 or a similar course.

Course Notes

Students are not allowed to use laptops or smartphones during class, unless used for class presentation. Homework assignments will be posted on Blackboard.

Required Books

- Salganik, M. J. (2017). *Bit by bit: social research in the digital age*. Princeton University Press
- Pennebaker, J. (2011). *The secret life of pronouns: What our words say about us*. New York, NY: Bloomsbury

Description and Assessment of Assignments

1. Paper presentation. Each student will present a set of papers related to one of the topics discussed in class.
2. Reaction paragraphs. Students are asked to write a short note, one or two paragraphs in length, about their reaction to the reading assignments of the week. These can be a quick summary of the material, comments about the subject area, or a critique of a particular theory or experiment. I will read these paragraphs carefully before each class, and will use them to guide the discussion in class. Simply reading the first page of a paper and writing a summary of it will not count as a reaction paragraph.
3. Class Project. This class is project oriented, and group-based. The goal of the project is for students to get experience working in interdisciplinary groups to tackle specific social scientific problems, and bring together theory from the social sciences and NLP techniques from computer science to tackle that problem. This will include a project proposal presentation, three project update presentations, final project presentation, and a report. For project proposals, students will present a problem and a data collection method and/or dataset for which they want to analyze. Each presentation should be about 10-15mins. The goal of the project update presentations is to inform the class about the state of the project and brainstorm with other students on how to solve the remaining issues. Each update presentation should be around 10 minutes. For the final project presentation, each student/group will give a 15-20min presentation on their project. Students are expected to spend at least 80 hours working on their final project. The project report will be around 20 pages.

Grading Policy

- 15% Participation
- 20% Paper Presentations
- 25% Reaction Paragraphs
- 15% Project Status Updates
- 10% Final Project Presentation
- 15% Final Project Write up

Assignment Submission Policy

All assignments are due on Thursdays at 10am. Assignments turned in any later than 10:10am will be considered late. Students will be allowed a total of four late days that can be used on the assignments. In exceptional circumstances, arrangements must be made in advance of the due date to obtain an extension. Once you have used up your four late days, one additional day late will result in a 25% reduction in the total score, two additional days late will yield a 50% reduction, and no credit will be given for three or more additional days late. Late days are in units of days, not hours, so using up part of a day uses up the whole day. The final project report, plus code used, will be due on the day of the final exam.

Schedule and weekly learning goals

The schedule is tentative and subject to change.

Week 01, 08/26: Introduction to Computational Social Sciences 1/2

- Lazer, D., Hargittai, E., Freelon, D., Gonzalez-Bailon, S., Munger, K., Ognyanova, K., and Radford, J. (2021). Meaningful measures of human society in the twenty-first century. *Nature*, 595(7866):189–196
- Hofman, J. M., Watts, D. J., Athey, S., Garip, F., Griffiths, T. L., Kleinberg, J., Margetts, H., Mullainathan, S., Salganik, M. J., Vazire, S., Vespignani, A., and Yarkoni, T. (2021). Integrating explanation and prediction in computational social science. *Nature*, 595(7866):181–188
- Wagner, C., Strohmaier, M., Olteanu, A., Kicinman, E., Contractor, N., and Eliassi-Rad, T. (2021). Measuring algorithmically infused societies. *Nature*, 595(7866):197–204
- Adjerid, I. and Kelley, K. (2018). Big data in psychology: A framework for research advancement. *American Psychologist*, 73(7):899

Week 02, 09/02: Introduction to Computational Social Sciences 2/2

- Salganik, M. J. (2017). *Bit by bit: social research in the digital age*. Princeton University Press (Chapters 1-3)
- Kennedy, B., Ashokkumar, A., Boyd, R. L., and Dehghani, M. (2022). Text analysis for psychology: Methods, principles, and practices. In Dehghani, M. and Boyd, R. L., editors, *Handbook of Language Analysis in Psychology*. Guilford Press, New York, NY

Week 03, 09/09: Dictionary Methods 1/2

- Pennebaker, J. (2011). *The secret life of pronouns: What our words say about us*. New York, NY: Bloomsbury

Week 04, 09/16: Dictionary Methods 2/2

- Back, M. D., Kufner, A. C., and Egloff, B. (2010). The emotional timeline of september 11, 2001. *Psychological Science*, 21(10):1417–1419
- Pury, C. L. (2011). Automation can lead to confounds in text analysis: Back, kufner, and egloff (2010) and the not-so-angry americans. *Psychological science*, 22(6):835
- Back, M. D., Kufner, A. C., and Egloff, B. (2011). “automatic or the people?”: Anger on september 11, 2001, and lessons learned for the analysis of large digital data sets. *Psychological Science*, 22(6):837
- Mehl, M. R., Raison, C. L., Pace, T. W., Arevalo, J. M., and Cole, S. W. (2017). Natural language indicators of differential gene regulation in the human immune system. *Proceedings of the National Academy of Sciences*, 114(47):12554–12559
- Choose two:
 - Iliev, R., Hoover, J., Dehghani, M., and Axelrod, R. (2016). Linguistic positivity in historical texts reflects dynamic environmental and psychological factors. *Proceedings of the National Academy of Sciences*, 113(49):E7871–E7879
 - Kaplan, D. M., Raison, C. L., Milek, A., Tackman, A. M., Pace, T. W., and Mehl, M. R. (2018). Dispositional mindfulness in daily life: A naturalistic observation study. *PloS one*, 13(11):e0206029
 - Boyd, R. L., Blackburn, K. G., and Pennebaker, J. W. (2020). The narrative arc: Revealing core narrative structures through text analysis. *Science advances*, 6(32):eaba2196

Week 05, 09/23: Differential Language Analysis & Project Proposals

- Park, G., Schwartz, H. A., Eichstaedt, J. C., Kern, M. L., Kosinski, M., Stillwell, D. J., Ungar, L. H., and Seligman, M. E. (2015). Automatic personality assessment through social media language. *Journal of personality and social psychology*, 108(6):934
- Eichstaedt, J. C., Schwartz, H. A., Kern, M. L., Park, G., Labarthe, D. R., Merchant, R. M., Jha, S., Agrawal, M., Dziurzynski, L. A., Sap, M., et al. (2015). Psychological language on twitter predicts county-level heart disease mortality. *Psychological science*, 26(2):159–169
- Brown, N. J. L. and Coyne, J. (2018). No evidence that twitter language reliably predicts heart disease: A reanalysis of eichstaedt et al. (2015a)
- Eichstaedt, j. C., Schwartz, H. A., Giorgi, S., Kern, M. L., Park, G., Sap, M., Labarthe, D. R., Larson, E. E., Seligman, M., and Ungar, L. H. (2018). More evidence that twitter language predicts heart disease: A response and replication
- Curtis, B., Giorgi, S., Buffone, A. E., Ungar, L. H., Ashford, R. D., Hemmons, J., Summers, D., Hamilton, C., and Schwartz, H. A. (2018). Can twitter be used to predict county excessive alcohol consumption rates? *PloS one*, 13(4):e0194290

Week 06, 09/30: Distributed Dictionary Representations

- Garten, J., Hoover, J., Johnson, K. M., Boghrati, R., Iskiwitch, C., and Dehghani, M. (2018). Dictionaries and distributions: Combining expert knowledge and large scale textual data content analysis. *Behavior research methods*, 50(1):344–361
- Dehghani, M., Johnson, K., Hoover, J., Sagi, E., Garten, J., Parmar, N. J., Vaisey, S., Iliev, R., and Graham, J. (2016). Purity homophily in social networks. *Journal of Experimental Psychology: General*, 145(3):366
- Hoover, J., Johnson, K., Boghrati, R., Graham, J., and Dehghani, M. (2018). Moral framing and charitable donation: Integrating exploratory social media analyses and confirmatory experimentation. *Collabra: Psychology*, 4(1)
- Bhatia, S. (2017). Associative judgment and vector space semantics. *Psychological Review*, 124(1):1
- Wang, S.-Y. N. and Inbar, Y. (2021a). Moral-language use by us political elites. *Psychological Science*, 32(1):14–26

Week 07, 10/07: Neural Networks 1/2

- Lin, Y., Hoover, J., Portillo-Wightman, G., Park, C., Dehghani, M., and Ji, H. (2018). Acquiring background knowledge to improve moral value prediction. In *2018 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, pages 552–559. IEEE
- Mooijman, M., Hoover, J., Lin, Y., Ji, H., and Dehghani, M. (2018). Moralization in social networks and the emergence of violence during protests. *Nature human behaviour*, 2(6):389
- Kennedy, B., Atari, M., Davani, A. M., Hoover, J., Omrani, A., Graham, J., and Dehghani, M. (2021). Moral concerns are differentially observable in language. *Cognition*, 212:104696
- Hoover, J., Atari, M., Davani, A. M., Kennedy, B., Portillo-Wightman, G., Yeh, L., Kogon, D., and Dehghani, M. (2019). Bound in hatred: The role of group-based morality in acts of hate
- Hoover, J., Atari, M., Mostafazadeh Davani, A., Kennedy, B., Portillo-Wightman, G., Yeh, L., and Dehghani, M. (2021). Investigating the role of group-based morality in extreme behavioral expressions of prejudice. *Nature Communications*, 12(1):4585
- Atari, M., Davani, A. M., Kogon, D., Kennedy, B., Saxena, N. A., Anderson, I., and Dehghani, M. (2021). Morally homogeneous networks and radicalism

Week 08, 10/14: Fall Recess

Week 09, 10/21: Neural Networks 2/2

- Sap, M., Horvitz, E., Choi, Y., Smith, N. A., and Pennebaker, J. (2020). Recollection versus imagination: Exploring human memory and cognition via neural language models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1970–1978
- Lucy, L., Demszky, D., Bromley, P., and Jurafsky, D. (2020). Content analysis of textbooks via natural language processing: Findings on gender, race, and ethnicity in texas us history textbooks. *AERA Open*, 6(3):2332858420940312
- Chakravarthula, S. N., Baucom, B. R., Narayanan, S., and Georgiou, P. (2021). An analysis of observation length requirements for machine understanding of human behaviors from spoken language. *Computer Speech & Language*, 66:101162
- Choose three:
 - Priniski, J. H., Mokhberian, N., Harandizadeh, B., Morstatter, F., Lerman, K., Lu, H., and Brantingham, P. J. (2021). Mapping moral valence of tweets following the killing of george floyd. *arXiv preprint arXiv:2104.09578*
 - Garten, J., Kennedy, B., Hoover, J., Sagae, K., and Dehghani, M. (2019). Incorporating demographic embeddings into language understanding. *Cognitive science*, 43(1):e12701
 - Forbes, M., Hwang, J. D., Shwartz, V., Sap, M., and Choi, Y. (2020). Social chemistry 101: Learning to reason about social and moral norms. *arXiv preprint arXiv:2011.00620*
 - Brady, W. J., McLoughlin, K., Doan, T. N., and Crockett, M. J. (2021). How social learning amplifies moral outrage expression in online social networks. *Science Advances*, 7(33)

Week 10, 10/28: Bias in NLP & Project Update I

- Caliskan, A., Bryson, J. J., and Narayanan, A. (2017). Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186
- Garg, N., Schiebinger, L., Jurafsky, D., and Zou, J. (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644
- Koenecke, A., Nam, A., Lake, E., Nudell, J., Quartey, M., Mengesha, Z., Toups, C., Rickford, J. R., Jurafsky, D., and Goel, S. (2020). Racial disparities in automated speech recognition. *Proceedings of the National Academy of Sciences*, 117(14):7684–7689
- Choose two:

- Charlesworth, T. E., Yang, V., Mann, T. C., Kurdi, B., and Banaji, M. R. (2021). Gender stereotypes in natural language: Word embeddings show robust consistency across child and adult language corpora of more than 65 million words. *Psychological Science*, 32(2):218–240
- Shah, D. S., Schwartz, H. A., and Hovy, D. (2020). Predictive biases in natural language processing models: A conceptual framework and overview. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5248–5264, Online. Association for Computational Linguistics
- Gonen, H. and Goldberg, Y. (2019). Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them. *arXiv preprint arXiv:1903.03862*

Week 11, 11/04: Other methods

- Voigt, R., Camp, N. P., Prabhakaran, V., Hamilton, W. L., Hetey, R. C., Griffiths, C. M., Jurgens, D., Jurafsky, D., and Eberhardt, J. L. (2017). Language from police body camera footage shows racial disparities in officer respect. *Proceedings of the National Academy of Sciences*, 114(25):6521–6526
- Camp, N. P., Voigt, R., Jurafsky, D., and Eberhardt, J. L. (2021). The thin blue waveform: Racial disparities in officer prosody undermine institutional trust in the police. *Journal of Personality and Social Psychology*
- Brady, W. J., Wills, J. A., Jost, J. T., Tucker, J. A., and Van Bavel, J. J. (2017). Emotion shapes the diffusion of moralized content in social networks. *Proceedings of the National Academy of Sciences*, 114(28):7313–7318
- Burton, J. W., Cruz, N., and Hahn, U. (2021). Reconsidering evidence of moral contagion in online social networks. *Nature Human Behaviour*, pages 1–7
- Wang, S.-Y. N. and Inbar, Y. (2021b). Re-examining the diffusion of moralized rhetoric from political elites: Effects of valence and ideology. *Under-review*
- Candia, C., Atari, M., Kteily, N., and Uzzi, B. (2021). Overuse of moral language dampens the diffusion of moralized content on social media. *Under-review*

Week 12, 11/11: Clinical & cognitive applications & Project Update II

- Resnik, P., Armstrong, W., Claudino, L., Nguyen, T., Nguyen, V.-A., and Boyd-Graber, J. (2015). Beyond lda: exploring supervised topic modeling for depression-related language in twitter. In *Proceedings of the 2nd Workshop on Computational Linguistics and Clinical Psychology: From Linguistic Signal to Clinical Reality*, pages 99–107

- Coppersmith, G., Leary, R., Crutchley, P., and Fine, A. (2018). Natural language processing of social media as screening for suicide risk. *Biomedical informatics insights*, 10:1178222618792860
- Mitchell, T. M., Shinkareva, S. V., Carlson, A., Chang, K.-M., Malave, V. L., Mason, R. A., and Just, M. A. (2008). Predicting human brain activity associated with the meanings of nouns. *science*, 320(5880):1191–1195
- Choose two:
 - Dehghani, M., Boghrati, R., Man, K., Hoover, J., Gimbel, S. I., Vaswani, A., Zevin, J. D., Immordino-Yang, M. H., Gordon, A. S., Damasio, A., et al. (2017). Decoding the neural representation of story meanings across languages. *Human brain mapping*, 38(12):6096–6106
 - Toneva, M. and Wehbe, L. (2019). Interpreting and improving natural-language processing (in machines) with natural language-processing (in the brain). *arXiv preprint arXiv:1905.11833*
 - Schwartz, D., Toneva, M., and Wehbe, L. (2019). Inducing brain-relevant bias in natural language processing models. *arXiv preprint arXiv:1911.03268*

Week 13, 11/18: Ethics

- Salganik, M. J. (2017). *Bit by bit: social research in the digital age*. Princeton University Press (Chapters 6-7)
- Alfano, M., Sullivan, E., and Ebrahimi Fard, A. (2022). Ethical pitfalls for natural language processing in psychology. In Dehghani, M. and Boyd, R. L., editors, *Handbook of Language Analysis in Psychology*. Guilford Press, New York, NY
- Skorburg, J. A. and Friesen, P. (2022). Ethical issues in text mining for mental health. In Dehghani, M. and Boyd, R. L., editors, *Handbook of Language Analysis in Psychology*. Guilford Press, New York, NY
- Hovy, D. and Spruit, S. L. (2016). The social impact of natural language processing. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 591–598
- Card, D., Henderson, P., Khandelwal, U., Jia, R., Mahowald, K., and Jurafsky, D. (2020). With little power comes great responsibility. *arXiv preprint arXiv:2010.06595*
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, pages 610–623, New York, NY, USA. Association for Computing Machinery
- Watch in class: Friends You Haven't Met Yet

Week 14, 11/25: Thanksgiving Holiday

Week 15, 12/02: Final project presentations

Statement on Academic Conduct and Support Systems

Academic Conduct

Plagiarism — presenting someone else’s ideas as your own, either verbatim or recast in your own words — is a serious academic offense with serious consequences. Please familiarize yourself with the discussion of plagiarism in *SCampus* in Section 11, *Behavior Violating University Standards* <https://scampus.usc.edu/1100-behavior-violating-university-standards-and-appropriate-sanctions/>. Other forms of academic dishonesty are equally unacceptable. See additional information in *SCampus* and university policies on scientific misconduct, <http://policy.usc.edu/scientific-misconduct/>.

Discrimination, sexual assault, and harassment are not tolerated by the university. You are encouraged to report any incidents to the *Office of Equity and Diversity* <http://equity.usc.edu/> or to the *Department of Public Safety* <http://capsnet.usc.edu/department/department-public-safety/online-forms/contact-us>. This is important for the safety whole USC community. Another member of the university community — such as a friend, classmate, advisor, or faculty member — can help initiate the report, or can initiate the report on behalf of another person. *The Center for Women and Men* <http://www.usc.edu/student-affairs/cwm/> provides 24/7 confidential support, and the sexual assault resource center webpage sarc@usc.edu describes reporting options and other resources.

Support Systems

A number of USC’s schools provide support for students who need help with scholarly writing. Check with your advisor or program staff to find out more. Students whose primary language is not English should check with the *American Language Institute* <http://dornsife.usc.edu/ali>, which sponsors courses and workshops specifically for international graduate students. *The Office of Disability Services and Programs* http://sait.usc.edu/academicssupport/centerprograms/dsp/home_index.html provides certification for students with disabilities and helps arrange the relevant accommodations. If an officially declared emergency makes travel to campus infeasible, *USC Emergency Information* <http://emergency.usc.edu/will> provide safety and other updates, including ways in which instruction will be continued by means of blackboard, teleconferencing, and other technology.

IMPORTANT: COVID-19 PROTOCOLS

Students must comply with all COVID-19 safety protocols outlined by federal, state, local, and university policies. These policies will likely evolve with the changing conditions of the COVID-19 pandemic and may include social distancing, the use of face coverings at all times, proof of vaccination, and regular COVID testing, among others. Depending on the policies outline by the above authorities, and the conditions of the class, the class might switch between meeting online and in person.