

**DSCI-531: Fairness in Artificial Intelligence**

**Units: 4**

**Spring 2024 MW 4:00-5:50pm**

**Room: OHE 132**

**Piazza:**

**<https://piazza.com/usc/spring2024/dsci531>**

**Instructor: Keith Burghardt**

**Office:** ISI 940

**Office Hours:** Immediately after class or by appointment

**Contact Info:** keithab@isi.edu, Zoom

**TA: Ashwin Rao**

**Office:** ISI 911 **OR** <https://usc.zoom.us/j/9026240639>

**Office Hours:** 10AM - 11AM PT Monday on Zoom or by appointment

**Contact Info:** mohanrao@usc.edu

**TA: Fiona Guo**

**Office:** ISI 921 **OR** <https://usc.zoom.us/my/siyiguo>

**Office Hours:** 10-11 Wednesday on Zoom or by appointment

**Contact Info:** siyiguo@isi.edu

**Catalog Description**

Our society's rapid algorithmification is fueled by data, but the reliance on data raises important questions. What are the latent biases hidden in the collected data? If that data was used to train machine learning algorithms, how did these biases impact predictions made by algorithms and systems that depend on them? Are the algorithmic decisions fair, or do they perpetuate stereotypes and fortify discrimination? As we come to rely on AI to make decisions in our lives and allow for synergistic relationship with technology, we need to build trust in AI by improving algorithmic fairness, accountability, transparency and explainability.

**Course Description**

The course will explore topics in the intersection of data, language, networks and algorithms with fairness and bias through quantitative analysis and hands-on exploration. The course will introduce students to basic and advanced fairness concepts, including methods and apply them to societal data to study fairness and bias, and to understand their effects on learning algorithms. While there are no formal prerequisites for the

class, students are expected to be proficient in Python and have working knowledge of algorithm design and data structures, and to have taken college level or above courses in linear algebra and statistics. AI and machine learning courses are strongly recommended.

## Learning Objectives

After completing the course, students will be able to do the following:

- Demonstrate a familiarity with the principles of ethical considerations for AI systems. This includes definitions of fairness in computer science and related literature.
- Understand and identify sources of bias in data analysis tasks, including text, image and tabular data.
- Understand the impacts of bias on machine learning tasks.
- Choose or design bias mitigation algorithms suitable for a particular task.
- Detect and assess biases in both datasets and trained machine learning models.
- Presenting technical results via technical reports.
- Generate thought provoking peer-reviews.
- Present your findings in the form of a short presentation.

## Recommended Preparation

Knowledge of Python 3; undergraduate coursework in linear algebra, basic probability and statistics; familiarity with artificial intelligence and machine learning at the level of DSCI 552.

## Course Notes

The course will be run as a lecture class with student participation strongly encouraged. There are weekly readings and students are encouraged to do the readings prior to the discussion in class. All of the course materials, including the readings, lecture slides, homeworks will be posted online on Blackboard. The class project is a significant aspect of this course and at the end of the semester, students will present their projects in class.

## Technological Proficiency and Hardware/Software Required

Students are expected to know how to program in a language such as Java, C++, or Python 3 (recommended). Students are also expected to have their own laptop or desktop computer where they can install and run software to do the homework assignments.

## Required Readings and Supplementary Materials

There are no required textbooks. The reading material is based on recently published technical papers available via the ACM/IEEE/Springer digital libraries. All USC students have automatic access to these digital archives. Lecture slides will be placed on Piazza and will be accessible to students before each lecture.

## Description and Assessment of Assignments

### Weekly Readings and Quizzes

There will be 5 random quizzes at the beginning of class. These will not be announced ahead of time, but will be announced in class. The questions will concern some basic ideas discussed in the class and/or the topics related to the recommended readings. Each quiz will be completed in class on a Monday. Your lowest score will be dropped.

### Homework Assignments

There will be four homework assignments designed to give students proficiency with understanding biases in different data modalities:

- **HW1: Ethics in the Research Process** - This homework will introduce students to the concepts of conducting fair and ethical research. The focus will be on historical understanding of the nature of IRB. Students will complete the CITI Human Subjects Research training as part of their assignment.
- **HW2: Bias in Data and Prediction** – Students will learn to apply basic data mining techniques to data. Students will design and conduct statistical tests on large datasets. These tests will be designed around fairness concepts, and how to leverage techniques to identify unfairness.

- **HW3: Bias in NLP** -- Students will study gender bias in named entity recognition. Solving this homework requires basic natural language processing techniques including transformer-based language models like BERT.
- **HW4: Bias in Networks** -- In this homework, students will learn and apply basic network techniques to uncover the presence of gender bias in a network. Are women more or less represented in the network? Do they tend to occupy positions of higher or lower centrality than men?

**Note:** In both written and programming assignments, the completeness and the clarity of your description and analysis will matter as much as the final correct answer. Sending just a single final value (even if correct) is not enough. See the table below:

Grade Component	Meets Expectations (75%-100%)	Approaches Expectations (50%-75%)	Needs Improvement (0%-50%)
<b>Completeness (50%)</b>	All parts of the question are addressed. If the task was to a) select a machine learning algorithm, b) train, and c) validate the model, all three parts are completed.	Most parts of the question are addressed. If the task was to a) select a machine learning algorithm, b) train, and c) validate the model, the student selected and trained the model, but the validation part is missing or is incomplete.	The main question is not addressed.  The answer is irrelevant to the task.  The analysis or evaluation of the issues and events is either vague or <i>completely</i> inaccurate.
<b>Clarity (25%)</b>	A non-expert (e.g., a fellow student) can understand the solutions. All concepts and used techniques are defined and explained. Whenever it is applicable, the solution is accompanied by illustrative plots that are explained and interpreted. Accompanied code is well commented and easy to follow.	The teacher (or other professional physicists) can understand the solution but a non-expert might have some trouble doing so. The solution has some minor shortcuts or some non-explained assumptions. Not every step of the analysis is explained, but it is still possible to follow the author's logic. The code is not well commented but it is still possible to follow it.	It is hard to follow the solutions. The solution has some major shortcuts and hidden assumptions. The analysis or evaluation of the issues and events is vague. The code is not well commented but it is either hard or impossible to follow it.
<b>Validity (25%)</b>	All calculations are correct. The final values are right. The interpretations and final conclusions are valid.	Small mistake in the code and/or calculations (e.g., a wrong sign, a missing constant). The final answer is close to the correct value (e.g., by a small factor; twice too large or twice too small; however, the general trend is correct).	Major mistakes in the code and/or in the analysis. The final values and conclusions are incorrect.

### **Course Project**

An integral part of this course is the course project, which builds on the topics and techniques covered in the class, focusing on extending and evaluating methods to solve problems. Students will write a written proposal for the project, carry out literature review, conduct the project, review two proposals by your peers, address the comments you received on your own peer-reviews, and then write a paper about the project, and present

the project in class. Students are encouraged to identify a new problem, apply or extend the methods they learned in class to propose an approach to solve the problem. Students will propose a novel project, do the research and build a proof-of-concept, write a report about the work, and present the work in class. Emphasis is placed on quantitative evaluation of the approach. Each student will team up with a classmate (2-3 students) to do an independent project based on the topics covered in the class.

The objective of this assignment is to a) explore literature regarding data science and machine learning, b) synthesize the acquired knowledge in the form of article, c) learn how to write peer-review comments, d) learn how to respond to peer-review comments, e) be able to summarize a weeks-long project in a form of a condensed, short presentation.

### Structure and Formatting:

We encourage you to use the LaTeX template <https://www.overleaf.com/read/crjbrffthg> that we prepared for you in Overleaf. If you use a WYSIWYG editor, please remember to submit your article in the PDF format (not docx, rtf or odt). In the paper, you must provide a link to a GitHub repository with the relevant code, scripts, or notebooks. Python 3.8+ is preferred, but in principle, you are free to use any language of your choice - as long as the code is clear and well commented (the reader should be able to go to your repository, clone your repository and run your code without getting any errors).

### Project Timeline:

- Week 1-2: Identify team members.
- Week 3-4: Choose your topic. Prepare and post a work plan/proposal **by Wednesday, January 31, at 4:00 p.m.**
- Week 6: Find relevant literature. Read about your topic. Prepare a literature review **by Wednesday, February 14, at 4:00 p.m.**
  - The literature review should be about 2-3 pages
  - It needs to be clear what your project's goal is, the methodology you plan to use, and how you plan to evaluate. For example, it's no longer enough to say you'll be looking at bias-what form of bias? In what data set? This needs to be very clear. Feel free to reuse text from your project proposal.
  - The goal of a literature review is to understand and outline what has already been done, to better inform your work. How was your data set created, and how has it already been used? What methods have been used to do what you're setting out to do? Have people already done what you want to do, and if it's close then what is the difference? Think of a "related works" section in a paper.
  - Both partners must submit the literature review on Blackboard. If you don't submit but your partner does, you will still get a score of 0 for this portion.
- Week 6-8: Make a plan for your article. Decide which aspects you are going to describe and which leave out. After all, you have limited space (only a couple of pages, including figures and bibliography). Complete the necessary coding and calculations. Prepare plots and figures.
- Week 9: Write the first version of your report. You should have a draft by **March 6**. Proofread your report. Make sure that all key terms are defined. Make sure that the report has the right structure (abstract, introduction, the main content, discussion/summary, and bibliography). Remember, that the list of references at the end of your report is not enough - your sources must be cited in the article.
- Week 11: Prepare a pdf of your article (this is your Midterm Report). Make sure that the number of words is below the maximum limit (3,500). Make sure that your name, affiliation, abstract and paper title are visible on the first page. Submit the pdf using Blackboard **by Wednesday, March 27, not later than 4:00 p.m.**

- Week 12: Choose two reports prepared by your peers. Read those reports. We will provide you your team pairs for the midterm reviews. Each team will look at the midterm reports from their two assigned groups. Your review will consist of the following:
  - Summary
  - 3 Strengths
  - 3 Weaknesses

Give each author suggestions on how they can improve the papers. To make sure that each team will receive an equal number of comments. You should complete this action **by Monday, April 8, at 4:00 p.m.** When uploading your reviews, please make 2 files (one for each team) and a filename equal to the last names of each team member with " \_ " between each name. **Make sure to upload both reviews on Blackbaord.**

- Read the suggestions that you received from your peers. Address them (either incorporate the suggested changes or challenge them, describing why you think those changes would not improve the quality of your article).
- Prepare an in-class presentation (length depends on the number of projects, but plan 5 minute presentation) **on Monday, April 22 and Wednesday, April 24 (in class).**
- Submit your final report **by Wednesday, May 3 at 4:00 p.m.**

*Sample project:*

“Ageism in Traffic Policing.” the goal of the project is to explore the nature of police stops in the USA, with a particular focus on bias. Using a large corpus of police stops, the students conducted an analysis of the features that are indicative of the outcome of a traffic stop. They found that ageism is rampant in traffic policing with older subjects more likely to receive lenient outcomes. Students found that some patterns are correlated with the political leaning of the state, with “red” states presenting less ageism than “blue” states.

**Grading Breakdown**

**Homework:** There will be 4 homework assignments.

**Quizzes:** While there are 5 quizzes in class, only the 4 best scores are used for grade calculation.

**Class Participation:** Students are expected to attend every class and actively participate in the discussion.

**Project:** Projects will be graded on novelty, technical soundness, and the quality of evaluation. Reports and presentations will be graded according to the project grading rubric and the quality and clarity of presentation.

Assignment	Points	% of Grade
Homeworks	20	20
Quizzes	30	30
Class participation	5	5
Project proposal	3	3
Literature review	5	5
Midterm report	8	8
Peer review	4	4

Project report	15	15
Project presentation	10	10
<b>TOTAL</b>	<b>100</b>	<b>100</b>

### **Assignment Submission Policy**

Homework assignments are due at 4pm on the due date and should be submitted on Blackboard. A total of 7 late days can be used. For each submission, a maximum of 2 late days are allowed. After that each day you will lose 20% of the possible points for the assignment. After one week, the assignment cannot be submitted. No late days are allowed for quizzes. Late submission of quizzes will receive 0 points.

### **Course Schedule: A Weekly Breakdown (next page)**

	Topics/Daily Activities	Readings and Homework <i>It is expected that students read these in advance of each week and come to class prepared to discuss them.</i>	Deliverable / Due Dates
<b>Week 1</b> Jan 8 & 10 Dr. Burghardt	<b>Course Introduction &amp; Research Ethics</b>	<ul style="list-style-type: none"> <li>Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., &amp; Galstyan, A. (2019). A survey on bias and fairness in machine learning. <i>arXiv preprint arXiv:1908.09635</i>.</li> <li>Olteanu, Castillo &amp; Kiciman "Social data: Biases, methodological pitfalls, and ethical boundaries."</li> <li>Angwin et al. Machine bias. Propublica (2016) <a href="https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing">https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing</a></li> </ul>	Homework 1 Assigned
<b>Week 2</b> Jan 17 Dr. Burghardt	<b>Unsupervised methods: clustering, PCA, etc</b>	<ul style="list-style-type: none"> <li>Gira, Nizar, Michel Crucianu, and Nozha Boujemaa. "Unsupervised and semi-supervised clustering: a brief survey." A review of machine learning techniques for processing multimedia content 1 (2004): 9-16.</li> </ul>	Homework 1 Due (1/17)
<b>Week 3</b> Jan 22 & 24 Dr. Burghardt	<b>Supervised methods: Regression, Explanation, etc.</b>	<ul style="list-style-type: none"> <li>Bodo Winter, Tutorial on linear models <a href="http://www.bodowinter.com/tutorial/bw_LME_tutorial1.pdf">http://www.bodowinter.com/tutorial/bw_LME_tutorial1.pdf</a></li> <li>Bodo Winter, Tutorial on mixed effects analysis, <a href="http://www.bodowinter.com/tutorial/bw_LME_tutorial2.pdf">http://www.bodowinter.com/tutorial/bw_LME_tutorial2.pdf</a></li> <li>Shapley</li> <li>Fennell, Zou &amp; Lerman (2019) "Explaining and Predicting Behavioral Data with Structured Feature Space Decomposition", in <i>EPJ Data Science</i></li> </ul>	Homework 2 Assigned  Quiz assigned Monday - in Class and Piazza

<p><b>Week 4</b> Jan 29 &amp; Jan 31 Dr. Burghardt</p>	<p><b>Definitions and fair machine learning</b></p>	<ul style="list-style-type: none"> <li>• Kleinberg, Mullainathan, and Raghavan. 2016. Inherent trade-offs in the fair determination of risk scores (2016) <a href="https://arxiv.org/abs/1609.05807">https://arxiv.org/abs/1609.05807</a></li> <li>• Pleiss, G., Raghavan, M., Wu, F., Kleinberg, J., &amp; Weinberger, K. Q. (2017). On fairness and calibration. In <i>Advances in Neural Information Processing Systems</i> (pp.5680-5689). <a href="http://papers.nips.cc/paper/7151-on-fairness-and-calibration">http://papers.nips.cc/paper/7151-on-fairness-and-calibration</a></li> <li>• Corbett-Davies, Sam, et al. "Algorithmic decision making and the cost of fairness." <i>Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i>. ACM, 2017. <a href="https://dl.acm.org/citation.cfm?id=3098095">https://dl.acm.org/citation.cfm?id=3098095</a></li> </ul>	<p>Project Proposals Due</p>
<p><b>Week 5</b> Feb 5 &amp; 7 Dr. Burghardt</p>	<p><b>Bias in Data</b></p>	<ul style="list-style-type: none"> <li>• Ruths &amp; Pfeffer "Social Media for Large Studies of Behavior"</li> <li>• Barocas, S., &amp; Selbst, A. D. (2016). Big data's disparate impact. <i>Calif. L. Rev.</i>, 104, 671. Part 1. <a href="http://www.cs.yale.edu/homes/jf/BarocasSelbst.pdf">http://www.cs.yale.edu/homes/jf/BarocasSelbst.pdf</a></li> <li>• Vaupel, J. W., &amp; Yashin, A. I. (1985). Heterogeneity's ruses: some surprising effects of selection on population dynamics. <i>The American Statistician</i>, 39(3), 176-185. <a href="https://www.jstor.org/stable/2683925?seq=1#metadata_info_tab_contents">https://www.jstor.org/stable/2683925?seq=1#metadata_info_tab_contents</a></li> <li>• Lerman, K. (2018). Computational social scientist beware: Simpson's paradox in behavioral data. <i>Journal of Computational Social Science</i>, 1(1), 49-58. <a href="https://arxiv.org/abs/1710.08615">https://arxiv.org/abs/1710.08615</a></li> <li>• Alipourfard, N., Fennell, P. G., &amp; Lerman, K. (2018). Using Simpson's Paradox to Discover Interesting Patterns in Behavioral Data. In <i>Twelfth International AAAI Conference on Web and Social Media</i>. <a href="https://link.springer.com/article/10.1007/s42001-017-0007-4">https://link.springer.com/article/10.1007/s42001-017-0007-4</a></li> </ul>	<p>Homework 2 due (2/7)</p> <p>Homework 3 assigned</p> <p>Quiz assigned Monday - in Class and Piazza</p>

<p><b>Week 6</b> Feb 12 &amp; 14 Dr. Burghardt</p>	<p><b>Image fairness and domain adaptation</b></p>	<ul style="list-style-type: none"> <li>● <a href="https://towardsdatascience.com/understanding-domain-adaptation-5baa723ac71f">https://towardsdatascience.com/understanding-domain-adaptation-5baa723ac71f</a></li> <li>● Hal Daumé III. 2007. <i>Frustratingly Easy Domain Adaptation</i>. In <i>Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics</i>, pages 256–263, Prague, Czech Republic. Association for Computational Linguistics.</li> <li>● Tzeng, E., Hoffman, J., Saenko, K., &amp; Darrell, T. (2017). Adversarial Discriminative Domain Adaptation. <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)</i>. <a href="https://openaccess.thecvf.com/content_cvpr_2017/html/Tzeng_Adversarial_Discriminative_Domain_CVPR_2017_paper.html">https://openaccess.thecvf.com/content_cvpr_2017/html/Tzeng_Adversarial_Discriminative_Domain_CVPR_2017_paper.html</a></li> <li>● Takiddin A, Schneider J, Yang Y, Abd-Alrazaq A, Househ M. Artificial Intelligence for Skin Cancer Detection: Scoping Review. <i>J Med Internet Res</i>. 2021 Nov 24;23(11):e22934. doi: 10.2196/22934. PMID: 34821566; PMCID: PMC8663507. <a href="https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8663507/">https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8663507/</a></li> <li>● Min Kyung Lee and Katherine Rich. 2021. Who Is Included in Human Perceptions of AI?: Trust and Perceived Fairness around Healthcare AI and Cultural Mistrust. In <i>Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems (CHI '21)</i>. Association for Computing Machinery, New York, NY, USA, Article 138, 1–14. <a href="https://doi.org/10.1145/3411764.3445570">https://doi.org/10.1145/3411764.3445570</a></li> </ul>	
------------------------------------------------------------	----------------------------------------------------	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--

<p><b>Week 7</b> Feb 21 Dr. Burghardt</p>	<p><b>Basics of NLP: Text mining and word embeddings</b></p>	<ul style="list-style-type: none"> <li>● Bolukbasi, Chang, Zou, Saligrama &amp; Kalai "Man is to computer programmer as woman is to homemaker? debiasing word embeddings"</li> <li>● Zhao, Wang, Yatskar, Ordonez &amp; Chang "Men also like shopping: Reducing gender bias amplification using corpus-level constraints"</li> <li>● Zhao, Wang, Yatskar, Ordonez &amp; Chan "Gender bias in coreference resolution: Evaluation and debiasing methods"</li> </ul>	<p>Lit review due</p>
---------------------------------------------------	----------------------------------------------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------

<p><b>Week 8</b> Feb 26 &amp; Feb 28 Dr. Burghardt</p>	<p><b>Bias and Fairness in NLP</b></p>	<ul style="list-style-type: none"> <li>● Pak, A., &amp; Paroubek, P. (2010, May). Twitter as a corpus for sentiment analysis and opinion mining. In LREc (Vol. 10, No. 2010, pp. 1320-1326). <a href="https://lexitron.nectec.or.th/public/LREC-2010_Malta/pdf/385_Paper.pdf">https://lexitron.nectec.or.th/public/LREC-2010_Malta/pdf/385_Paper.pdf</a></li> <li>● Arvind Narayanan</li> <li>● S. Golder and M. Macy, Diurnal and Seasonal Mood Vary with Work, Sleep, and Daylength Across Diverse Cultures, Science Vol. 333 no. 6051 pp. 1878-1881, 2011. <a href="https://science.sciencemag.org/content/333/6051/1878/">https://science.sciencemag.org/content/333/6051/1878/</a></li> <li>● <b>More... Language reveals</b></li> <li>● Mehrabi, Ninareh, et al. "Man is to person as woman is to location: Measuring gender bias in named entity recognition." <i>Proceedings of the 31st ACM Conference on Hypertext and Social Media</i>. 2020.</li> <li>● Zhao, Jieyu, et al. "Men also like shopping: Reducing gender bias amplification using corpus-level constraints." arXiv preprint arXiv:1707.09457 (2017).</li> <li>● Feng et al., From Pretraining Data to Language Models to Downstream Tasks:</li> <li>● Tracking the Trails of Political Biases Leading to Unfair NLP Models (<a href="https://aclanthology.org/2023.acl-long.656.pdf">https://aclanthology.org/2023.acl-long.656.pdf</a>)</li> <li>● Santurkar et al., Whose Opinions Do Language Models Reflect? (<a href="https://arxiv.org/pdf/2303.17548.pdf">https://arxiv.org/pdf/2303.17548.pdf</a>)</li> <li>● Wan et al., "Kelly is a Warm Person, Joseph is a Role Model": Gender Biases in LLM-Generated Reference Letters ( <a href="https://arxiv.org/pdf/2310.09219.pdf">https://arxiv.org/pdf/2310.09219.pdf</a>)</li> <li>● Dhamala et al., BOLD: Dataset and Metrics for Measuring Biases in Open-Ended Language Generation (<a href="https://arxiv.org/pdf/2101.11718.pdf">https://arxiv.org/pdf/2101.11718.pdf</a>)</li> <li>● Lucy et al., Gender and Representation Bias in GPT-3 Generated Stories (<a href="https://aclanthology.org/2021.nuse-1.5.pdf">https://aclanthology.org/2021.nuse-1.5.pdf</a>)</li> <li>● Salinas et al., The Unequal Opportunities of Large Language Models: Revealing Demographic Bias through Job</li> </ul>	<p>Homework 3 due (2/28) Homework 4 assigned</p>
----------------------------------------------------------------	----------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------------------------------------------------

		<p>Recommendations  <a href="https://arxiv.org/pdf/2308.02053.pdf">https://arxiv.org/pdf/2308.02053.pdf</a></p> <ul style="list-style-type: none"> <li>Cheng et al., Marked Personas: Using Natural Language Prompts to Measure Stereotypes in Language Models  <a href="https://aclanthology.org/2023.acl-long.84.pdf">https://aclanthology.org/2023.acl-long.84.pdf</a></li> </ul>	
<p><b>Week 9</b>  Mar 4 &amp; 6  Dr. Burghardt</p>	<p><b>Basics of Networks</b></p>	<ul style="list-style-type: none"> <li>Easley &amp; Kleinberg "Networks, crowds, and markets" [Chapter 2]</li> <li>Fortunato &amp; Hric "Community detection in networks: A user guide"</li> <li>Mehrabi, Morstatter, Peng &amp; Galstyan "Debiasing Community Detection: The Importance of Lowly-Connected Nodes"</li> </ul>	<p>Quiz assigned Monday - in Class and Piazza</p>
<p><b>Mar 11,13</b></p>		<p><b>SPRING BREAK</b></p>	
<p><b>Week 10</b>  Mar 18 &amp; 20  Dr. Burghardt</p>	<p><b>Fairness &amp; Bias in Networks</b></p>	<ul style="list-style-type: none"> <li>Kooti, F., Hodas, N. O., &amp; Lerman, K. (2014). Network weirdness: Exploring the origins of network paradoxes. In <i>Eighth International AAAI Conference on Weblogs and Social Media</i>. <a href="https://www.aaai.org/ocs/index.php/ICWSM/ICWSM14/paper/viewFile/8101/8127">https://www.aaai.org/ocs/index.php/ICWSM/ICWSM14/paper/viewFile/8101/8127</a></li> <li>Lee, E., Karimi, F., Wagner, C., Jo, H. H., Strohmaier, M., &amp; Galesic, M. (2019). Homophily and minority-group size explain perception biases in social networks. <i>Nature human behaviour</i>, 1-10. <a href="https://www.nature.com/articles/s41562-019-0677-4">https://www.nature.com/articles/s41562-019-0677-4</a></li> <li>D Liben-Nowell &amp; J Kleinberg, "The link prediction problem for social networks." <i>Journal of the American Society for Information Science and Technology</i>, Vol. 58, No. 7. (May 2007), pp. 1019-1031.</li> <li>Stoica, AA, Riederer, C and Chaintreau, A. 2018. Algorithmic Glass Ceiling in Social Networks: The effects of social recommendations on network diversity. In <i>Proceedings of the 2018 World Wide Web Conference (WWW '18)</i>. 923-932. DOI: <a href="https://doi.org/10.1145/3178876.3186140">https://doi.org/10.1145/3178876.3186140</a></li> </ul>	<p>Homework 4 due (3/20)</p>
<p><b>Week 11</b>  Mar 25 &amp; 27</p>	<p><b>Debiasing Methods</b></p>	<ul style="list-style-type: none"> <li>Kamiran, Faisal, and Toon Calders. "Data preprocessing techniques for</li> </ul>	

Dr. Burghardt		<p>classification without discrimination." Knowledge and Information Systems 33.1 (2012): 1-33.</p> <ul style="list-style-type: none"> <li>● Zhang et al. "Mitigating unwanted biases with adversarial learning." Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society. 2018.</li> <li>● Adversarial approaches</li> <li>● Fair representation of graphs DeepWalk, FairWalk,</li> <li>● Fair representation of NLP, etc. BERT, Fair Word embedding</li> </ul>	<p>Midterm report due</p> <p>Quiz assigned Monday - in Class and Piazza</p>
---------------	--	-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------

<p><b>Week 12</b> Apr 1 &amp; 3 Dr. Burghardt</p>	<p><b>Privacy &amp; Fraud</b></p>	<ul style="list-style-type: none"> <li>• Ferrara, Varol, Davis, Menczer &amp; Flammini "The rise of social bots".</li> <li>• Pfeffer, Mayer &amp; Morstatter "Tampering with Twitter's Sample API".</li> <li>• Kosinski, M., Stillwell, D., and Graepel, T. (2013). Private traits and attributes are predictable from digital records of human behavior. <i>Proceedings of the National Academy of Sciences</i>, 110(15):5802-5805</li> <li>• Zheleva E., Getoor L. (2011) Privacy in Social Networks: A Survey. In: Aggarwal C. (eds) <i>Social Network Data Analytics</i>. Springer, Boston, MA <a href="https://link.springer.com/chapter/10.1007/978-1-4419-8462-3_10">https://link.springer.com/chapter/10.1007/978-1-4419-8462-3_10</a></li> </ul>	<p>Peer reviews due (4/3)</p>
<p><b>Week 13</b> Apr 8 &amp; 10 Dr. Burghardt</p>	<p><b>Algorithmic Bias and Feedback</b></p>	<ul style="list-style-type: none"> <li>• Salganik, M. J., Dodds, P. S., &amp; Watts, D. J. (2006). Experimental study of inequality and unpredictability in an artificial cultural market. <i>science</i>, 311(5762), 854-856. <a href="https://science.sciencemag.org/content/311/5762/854">https://science.sciencemag.org/content/311/5762/854</a></li> <li>• Ricardo Baeza-Yates "Bias on the Web" <i>Communications of the ACM</i>, June 2018, Vol. 61 No. 6, Pages 54-61 <a href="https://www.researchgate.net/profile/Ricardo_Baeza-Yates/publication/325330277_Bias_on_the_web/links/5b1576440f7e9bda0ffcc999/Bias-on-the-web.pdf">https://www.researchgate.net/profile/Ricardo_Baeza-Yates/publication/325330277_Bias_on_the_web/links/5b1576440f7e9bda0ffcc999/Bias-on-the-web.pdf</a></li> <li>• Lerman, K., &amp; Hogg, T. (2014). Leveraging position bias to improve peer recommendation. <i>PloS one</i>, 9(6), e98914. <a href="https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0098914">https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0098914</a></li> </ul>	<p>Quiz assigned Monday - in Class and Piazza</p>

<b>Week 14</b> Apr 15 & 17 Dr. Burghardt	<b>Researcher Bias &amp; Bonus Topic</b>	<ul style="list-style-type: none"> <li>• MacCoun, R., &amp; Perlmutter, S. (2015). Blind analysis: hide results to seek the truth. <i>Nature News</i>, 526(7572), 187. <a href="https://www.nature.com/news/blind-analysis-hide-results-to-seek-the-truth-1.18510">https://www.nature.com/news/blind-analysis-hide-results-to-seek-the-truth-1.18510</a></li> <li>• Silberzahn, R., Uhlmann, E. L., Martin, D. P., Anselmi, P., Aust, F., Awtrey, E. C., ... Nosek, B. A. (2017, April 24). Many analysts, one dataset: Making transparent how variations in analytical choices affect results. <a href="https://doi.org/10.1177/2515245917747646">https://doi.org/10.1177/2515245917747646</a></li> <li>• Arvind Narayanan’s talk “How to recognize AI snake oil” <a href="https://www.cs.princeton.edu/~arvindn/talks/MIT-STS-AI-snakeoil.pdf">https://www.cs.princeton.edu/~arvindn/talks/MIT-STS-AI-snakeoil.pdf</a></li> </ul>	
<b>Week 15</b> Dr. Burghardt Apr 22 & 24	<b>Project Presentations</b>		Project Presentation slides due
<b>FINAL</b> May 3	<b>Project Final Report</b>		Project Final Report due

## **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 Part B, Section 11, “Behavior Violating University Standards” [policy.usc.edu/scampus-part-b](http://policy.usc.edu/scampus-part-b). Other forms of academic dishonesty are equally unacceptable. See additional information in SCampus and university policies on scientific misconduct, [policy.usc.edu/scientific-misconduct](http://policy.usc.edu/scientific-misconduct).

### **Support Systems:**

*Counseling and Mental Health - (213) 740-9355 – 24/7 on call*  
[studenthealth.usc.edu/counseling](http://studenthealth.usc.edu/counseling)

Free and confidential mental health treatment for students, including short-term psychotherapy, group counseling, stress fitness workshops, and crisis intervention.

*National Suicide Prevention Lifeline - 1 (800) 273-8255 – 24/7 on call*  
[suicidepreventionlifeline.org](http://suicidepreventionlifeline.org)

Free and confidential emotional support to people in suicidal crisis or emotional distress 24 hours a day, 7 days a week.

*Relationship and Sexual Violence Prevention and Services (RSVP) - (213) 740-9355(WELL), press "0" after hours – 24/7 on call*

[studenthealth.usc.edu/sexual-assault](http://studenthealth.usc.edu/sexual-assault)

Free and confidential therapy services, workshops, and training for situations related to gender-based harm.

*Office of Equity and Diversity (OED)- (213) 740-5086 | Title IX – (213) 821-8298*

[equity.usc.edu](http://equity.usc.edu), [titleix.usc.edu](http://titleix.usc.edu)

Information about how to get help or help someone affected by harassment or discrimination, rights of protected classes, reporting options, and additional resources for students, faculty, staff, visitors, and applicants. The university prohibits discrimination or harassment based on the following *protected characteristics*: race, color, national origin, ancestry, religion, sex, gender, gender identity, gender expression, sexual orientation, age, physical disability, medical condition, mental disability, marital status, pregnancy, veteran status, genetic information, and any other characteristic which may be specified in applicable laws and governmental regulations. The university also prohibits sexual assault, non-consensual sexual contact, sexual misconduct, intimate partner violence, stalking, malicious dissuasion, retaliation, and violation of interim measures.

*Reporting Incidents of Bias or Harassment - (213) 740-5086 or (213) 821-8298*

[usc-advocate.symplicity.com/care\\_report](http://usc-advocate.symplicity.com/care_report)

Avenue to report incidents of bias, hate crimes, and microaggressions to the Office of Equity and Diversity | Title IX for appropriate investigation, supportive measures, and response.

*The Office of Disability Services and Programs - (213) 740-0776*

[dsp.usc.edu](http://dsp.usc.edu)

Support and accommodations for students with disabilities. Services include assistance in providing readers/notetakers/interpreters, special accommodations for test taking needs, assistance with architectural barriers, assistive technology, and support for individual needs.

*USC Support and Advocacy - (213) 821-4710*

[uscса.usc.edu](http://uscса.usc.edu)

Assists students and families in resolving complex personal, financial, and academic issues adversely affecting their success as a student.

*Diversity at USC - (213) 740-2101*

[diversity.usc.edu](http://diversity.usc.edu)

Information on events, programs and training, the Provost's Diversity and Inclusion Council, Diversity Liaisons for each academic school, chronology, participation, and various resources for students.

*USC Emergency - UPC: (213) 740-4321, HSC: (323) 442-1000 – 24/7 on call*

[dps.usc.edu](http://dps.usc.edu), [emergency.usc.edu](http://emergency.usc.edu)

Emergency assistance and avenue to report a crime. Latest updates regarding safety, including ways in which instruction will be continued if an officially declared emergency makes travel to campus infeasible.

*USC Department of Public Safety - UPC: (213) 740-6000, HSC: (323) 442-120 – 24/7 on call*

[dps.usc.edu](http://dps.usc.edu)

Non-emergency assistance or information.