

INFO/CS 4300: Language and Information, Spring 2017

Course Syllabus

- **Time and place** TuTh 2:55pm-4:10pm, Hollister Hall B14
- **Instructor:** [Prof. Cristian Danescu-Niculescu-Mizil](#)
- **PhD TAs:** [Tom Davidson](#), [Xilun Chen](#)
- **Undergrad TAs:** [Ilan Filonenko](#), [Teddy Heidmann](#), [Eyvind Niklasson](#), [Pujaa Rajan](#), [Abigail Shchur](#), [Mukund Sudarshan](#), [Andrew Wang](#)
- **Course Piazza page** <https://piazza.com/cornell/spring2017/csinfo4300/home> (access code will be provided via CMS)
- **Course homepage** <http://www.cs.cornell.edu/Courses/cs4300/2017sp/>
- **Office hours:** Schedule listed on Piazza
- **Summary** How to make sense of the vast amounts of information available online, and how to relate it and to the social context in which it appears? This course introduces basic tools for retrieving and analyzing unstructured textual information from the web and social media. Applications include information retrieval (with human feedback), sentiment analysis and social analysis of text. The coursework will include programming projects that play on the interaction between knowledge and social factors.
- **Prerequisites:**
 - Linear algebra: strong performance in MATH 2940 or equivalent (or strong performance in INFO 2950);
 - Discrete math: strong performance in CS 2800 or equivalent (or strong performance in INFO 2950);
 - Programming proficiency: CS 2110 or equivalent;
 - Strong Python skills and familiarity with IPython Notebooks.
- **Related courses offered this semester at Cornell:**
 - [CS 4780 Machine Learning for Intelligent Systems](#)
 - [CS 5740 Natural Language Processing](#)
 - [CS 4744, LING 4424 Computational Linguistics](#)

Academic Integrity

We will strictly follow Cornell University's policies on academic integrity as outlined in the [Academic Integrity Handbook](#).

Any work submitted by a student in this course for academic credit will be the student's own work. For this course, collaboration is allowed only when it is made explicit in the assignment or project description. In case of doubt, contact the instructor.

All course materials are intellectual property belonging to the author. Students are not permitted to buy, sell or distribute any course materials without the express permission of the instructor. Such unauthorized behavior constitutes academic misconduct.

Late submissions and attendance

Attendance is mandatory, as for most lectures there will be no lecture slides. If you must miss a class, please email the instructor beforehand to provide an explanation. Late submissions will not be accepted, save for major medical or family events.

Electronic device policy

Notes for this class should be taken on paper. Use of electronic devices such as laptops and tablets will not be permitted during class (with the exception of specific activities). We are not plain evil, we are just following extensive research on the negative effects of in-class laptop use on learning.

Grading (subject to change)

Grades will be based on:

- participation (in-class or on Piazza) [10%];
- assignments/homeworks/quizzes [40%];
- midterm [20%];
- open ended final project [30%];

No auditing is allowed.

Midterm (subject to change)

The midterm will be administered in-class, likely the week before Spring Break or the week before that. Since this will be an in-class midterm there will be no makeup, so plan to attend.

SONA Credits

You can get extra credits for participating in experiments and research studies through [Science Research Participation System](#). You will receive 0.5% extra credit for each 30 minute study (or equivalent), up to a maximum of 1%.

Textbooks

- Manning, Raghavan, and Schütze. 2008. Introduction to Information Retrieval. Cambridge University Press.
- Jurafsky and Martin. 2009. Speech and Language Processing (2nd Edition). Pearson.

Course outline

The schedule and list of topics will be in **flux**. Here is a tentative outline:

| Week | Content |
|-------|---|
| 1 Th | Intro, Dimensions of information systems, Conversational behavior |
| 2 | Types and tokens, Document similarity |
| 3 | Vector space models, TF-IDF weighting |
| 4 | Indexing, Boolean search |
| 5 Th | Evaluation of IR systems |
| 6 | Ranked retrieval |
| 7 | Relevance feedback |
| 8 | Text classification |
| 9 | Midterm, rundown of textual features |
| 10 | Practical unsupervised text classification |
| 11 | Spring Break |
| 12 | Social features, Page Rank |
| 13 | Hubs and authorities Spectral analysis |
| 14 | Opinion mining, Trust, Deception |
| 15 | Final project presentations |
| 16 Tu | TBD |
