

CMU 94-775 UNSTRUCTURED DATA ANALYTICS FOR POLICY
(SPRING 2021 MINI-4 SECTION A4, 6 UNITS)

Instructor: George H. Chen (email: georgechen ♣ cmu.edu) — replace “♣” with an “at” symbol

Time and location: (lectures) Tuesdays and Thursdays, 1:30pm-2:50pm HBH 1206, (recitations) Fridays 4:50pm-6:10pm HBH 1206

TAs:

- Jingbo Jiang (jingboj ♣ andrew.cmu.edu)
- Xuejian Wang (xuejianw ♣ andrew.cmu.edu)

Office hours: TBA

Course webpage: www.andrew.cmu.edu/user/georgech/94-775/

Course description: Companies, governments, and other organizations now collect massive amounts of data such as text, images, audio, and video. How do we turn this heterogeneous mess of data into actionable insights? A common problem is that we often do not know what structure underlies the data ahead of time, hence the data often being referred to as “unstructured”. This course takes a practical approach to unstructured data analysis via a two-step approach:

- (1) We first examine how to identify possible structure present in the data via visualization and other exploratory methods.
- (2) Once we have clues for what structure is present in the data, we turn toward exploiting this structure to make predictions.

Many examples are given for how these methods help solve real problems faced by organizations. There is a final project in this course which must address a policy question.

How this course differs from 95-865 “Unstructured Data Analysis”: 95-865 emphasizes more of the technical skill development (assessed through two in-class exams involving coding), and does not have any sort of policy focus. 94-775 has a policy-focused final project instead of a final exam. 94-775 does not require cloud computing (part of 95-865 requires the use of Google Colab). Despite these differences, there is heavy material overlap between 94-775 and 95-865.

Learning objectives: By the end of the course, students are expected to have developed the following skills:

- Recall and discuss common methods for exploratory and predictive analysis of unstructured data
- Write Python code for exploratory and predictive data analysis
- Apply unstructured data analysis techniques discussed in class to solve problems faced by governments and companies

Skills are assessed by homework assignments, a quiz, and a final project. The homework and quiz are meant to be done individually whereas the final project is done in groups.

Prerequisites: If you are a Heinz student, then you must have either (1) passed the Heinz Python exemption exam, or (2) taken 90-819 “Intermediate Programming with Python” or 95-888 “Data-Focused Python”. If you are not a Heinz student and would like to take the course, please contact the instructor and clearly state what Python courses you have taken/what Python experience you have.

Instructional materials: There is no official textbook for the course. We will provide reading material as needed.

Homework: There are 3 homework assignments that give hands-on experience with techniques discussed in class. Assignments involve coding in Python and are submitted via Canvas.

Grading: Grades will be determined using the following weights:

| Assignment | Percentage of grade |
|--------------------------------|---------------------|
| HW1 | 8% |
| HW2 | 8% |
| HW3 (shorter than HW1 and HW2) | 4% |
| Final project proposal | 10% |
| Quiz | 35% |
| Final project | 35% |

Letter grades are assigned using a curve. Note that HW3 is designed to be shorter than HW1 and HW2 so that you have more time to work on the final project.

Cheating and plagiarism: In short, the only part of this course that is meant to be a group effort is the final project. While you are welcome to discuss homework problems with classmates, you must write up solutions to homework assignments on your own. At no time during the course should you have access to anyone else's code to any of the homework assignments including shared via instant messaging, email, Box, Dropbox, GitHub, Bitbucket, Amazon Web Services, etc. If part of your homework code or solutions uses an existing result (e.g., from a book, online resources), please cite your sources. For the quiz, your answers must reflect your work alone. Penalties for cheating range from receiving a 0 on an assignment to failing the course. In extreme circumstances, the instructor may file a case against you recommending the termination of your CMU enrollment.

Additional course policies:

Late homework: You are allotted a total of two late days that you may use however you wish for the homework assignments. By using a late day, you get a 24-hour extension without penalty. For example:

- You could use the two late days on two different assignments to turn them in each 1 day (24 hours) late without penalty.
- You could use both late days on a single homework assignment to turn it in 2 days (48 hours) late without penalty.

Note that you do *not* get fractional late days, e.g., you cannot use 1/2 of a late day to get a 12-hour extension. We will keep track of how many late days you have left based on the submission times of your homework assignments on Canvas. Once you have exhausted your late days, work you submit late will not be accepted. This policy only applies to homework; the quiz and final project must be submitted on time to receive any credit.

Re-grade policy: If you want an assignment regraded, please write up a note detailing your request and submit it to the instructor. Note that the entire assignment will be regraded and it is possible that your score may be lowered. The course staff will make it clear by what date re-grades for a particular assignment are accepted until. Re-grade requests submitted late will not be processed.

Course schedule (subject to revision; see course webpage for most up-to-date calendar): The course is roughly split into two parts. The first part is on exploratory data analysis in which given a dataset, we compute and visualize various aspects of it to try to understand its structure. The second part of the course turns toward making predictions once we have some idea of what structure underlies the data.

- Lecture 1 (Mar 23): Course overview
- Part I: Exploratory data analysis
 - Lecture 2 (Mar 25): Basic text analysis demo, co-occurrence analysis
 - Lecture 3 (Mar 30): Finding possibly related entities
 - Lecture 4 (Apr 1): Visualizing high-dimensional data with PCA
 - Lecture 5 (Apr 2): Manifold learning
 - Lecture 6 (Apr 6): Clustering I
 - **HW1 due Apr 6, 11:59pm**
 - Lecture 7 (Apr 8): Clustering II
 - Lecture 8 (Apr 9): Clustering III
 - Lecture 9 (Apr 13): Topic modeling with latent Dirichlet allocation
 - **HW2 and final project proposal due Apr 20, 11:59pm**
 - **Quiz on Apr 22 during regular lecture slot**
- Part II: Predictive data analysis:
 - Lecture 10 (Apr 23): Introduction to predictive data analytics
 - Lecture 11 (Apr 27): Introduction to neural nets and deep learning
 - Lecture 12 (Apr 29): Image analysis with convolutional neural nets
 - **HW3 due May 3, 11:59pm**
 - Lecture 13 (May 4): Time series analysis with recurrent neural nets
 - Lecture 14 (May 6): Course wrap-up
- **Final project presentations will be on Thursday May 6 (during part of lecture time) and Friday May 7 (during the recitation slot)**