

IS-392 Information Retrieval

Instructor: Dr. Christopher Markson, PhD

Class Location: Online (Asynchronous)

Email: crm23@njit.edu

Office Hours: by appointment via WebEx

Course Description: Web mining aims to discover useful information and knowledge from the Web hyperlink structure, page contents and usage logs. It has direct applications in e-commerce, Web analytics, information retrieval/filtering, personalization, and recommender systems. Employees knowledgeable about Web mining techniques and their applications are highly sought by major Web companies such as Google, Amazon, Yahoo, MSN and others who need to understand user behavior and utilize discovered patterns from terabytes of user profile data to design more intelligent applications. The primary focus of this course is on Web usage mining and its applications to business intelligence and biomedical domains. We learn techniques from machine learning, data mining, text mining, and databases to extract useful knowledge from the Web and other unstructured/semi-structured, hyper-textual, distributed information repositories. This data could be used for site management, automatic personalization, recommendation, and user profiling. Topics covered include crawling, indexing, ranking and filtering algorithms using text and link analysis, applications to search, classification, tracking, monitoring, and Web intelligence. Programming assignments give hands-on experience. A group project highlights class topics.

Required Background:

1. Knowledge of coding and data structures
2. Basic knowledge of database design and programming

Course Website: Canvas

Textbook:

Search Engines: Information Retrieval in Practice

Free download: <https://ciir.cs.umass.edu/irbook/>

Canvas: Additional material and resources will be found on the class website on Canvas. It will be modified and updated as the course progresses and will contain the most recent information.

Special Arrangement: There is an increasing demand from employers for the graduating students with the transferrable learning skills, including self-regulation, resilient, openness, communication skills, etc. These skills enable learners to continually upgrade the knowledge and all the learning related skills through their own self-motivated learning, which are more likely to happen in the working environment. However, improving these transferrable learning skills requires tremendous extra time and energy. In order to both encourage the learning activities which improve the transferrable learning skills and satisfy the need for just getting the credit to graduate, the full score (100%) is solely determined by the traditional learning through lecture for basic knowledge and the extra credit is determined by the unconventional learning through problem-based learning for advanced knowledge. To get the full score related to the lecture part, there is no requirement for the programming skill. All the related topics can be implemented in R, a popular data mining platform.

Credit: 3

Grade: Final grades will be based on:

Assignments: 25%

Midterm: 25%

Final Exam: 25%

Project 15%

Participation: 10%

The final letter grades for the semester are based solely on the points you earn (no curve).

Grade Points

A 90+

B+ 86-89

B 80-85

C+ 76-79

C 70-75

F 0-69

Lecture Schedule: The following is a tentative schedule and subject to change. Refer to the class web page for most recent information. All of the readings are from the main textbook for the course. For most topics, there is a laboratory part to apply the related algorithms to the given sample dataset to examine the output in R.

POLICIES:

Assignments (Homework and Project)

Homework will be submitted via Canvas electronically. Late homework will be penalized 10% of the available points (and another 10% will be deducted for every 24-hour period after the original due date). After two days beyond the deadline, I will no longer accept homework submissions (No exceptions).

Makeup Tests

Requests for makeup tests must be made in advance with the instructor and will only be approved if the reason is beyond your control.

Project

The project will consist of a report and Power point slides. This project should be more complicated than the homework assignments. The project should include the use of a web-based dataset, the analysis of the data, and code written in R or Python.

Academic Integrity Policy

The NJIT academic honor code is located at: <http://integrity.njit.edu/index.html>[Links to an external site.](#). This honor code applies in its entirety to this class. Violations will not be tolerated. In addition, students should familiarize themselves with NJIT's "Best Practices related to Academic Integrity" which is developed and published on the Provost's website (on the policies page).

TURNITIN Policy

NJIT uses Turnitin.com, a service that helps prevent plagiarism on student papers. I will be using the

Turnitin.com service at my discretion to determine the originality of student papers. If I submit your paper to Turnitin.com, it will be stored by Turnitin.com in their database as long as their service remains in

existence. If you object to this storage of your paper, you must let me know no later than two weeks after

the start of this class. If you object to the storage of your paper on Turnitin.com, I will utilize other services and techniques to check your work for plagiarism.

Disabilities

If you have a disability that may require some modification of seating, testing, or any other class

requirement; please let the Professor know so that appropriate arrangements can be made. Similarly let the Professor know if you have any emergency medical information about which to be aware, or if you need special arrangements in the event of building evacuation. See the Professor after class hours or schedule an appointment. Assistance is available from the Office of Student Disability Services (205

Campbell Hall; 973-596-3420). Be sure and fill out appropriate paperwork with this office during the first week of class.