

# CSE 435/535 Fall 2022 – Draft Syllabus

## Information Retrieval

Reg # 18180/18181

Lecture: Monday, Wednesday, Friday 1:00 pm – 1:50 pm (Buffalo time)

Knox 109

Instructor: Sougata Saha

### Description:

This course will introduce students to text-based information retrieval (IR) techniques, i.e. search engines. The course begins with the fundamentals of processing large-scale, multilingual text document collections. Various IR models such as the Boolean model, vector space model, and probabilistic models will be studied. Efficient indexing techniques for (i) general document collections, (ii) specialized collections (e.g. Wikipedia, biomedical, patents) and (iii) high velocity data such as social media will be discussed. Techniques for improving search efficiency, improving performance as well as evaluation methodology will be covered. The latter part of the course will focus on web search including link analysis techniques such as PageRank and HITS. The use of word vectors (Word2vec, GloVe) generated through neural models and their use in IR systems will be introduced. Students will work on programming projects (implemented on the AWS cloud computing platform) to gain hands-on expertise in building IR systems. This course provides the foundation for the follow-on course (CSE 635) which discusses natural language processing (NLP) and deeper text mining solutions.

**Prerequisites:** Programming expertise in Python, Linear Algebra

**Textbook:** Introduction to Information Retrieval by C. Manning, P. Raghavan, and H. Schütze, Cambridge University Press (2008, online version 2012)

Note: an online version of this book is available at <http://informationretrieval.org>

Other, more recent reference material will be made available on the piazza site during the semester.

**Instructor:** Sougata Saha, Dept. of Computer Science & Eng

email: [sougatas@buffalo.edu](mailto:sougatas@buffalo.edu)

office hours: TBA

### TAs:

Souvik Das. ([souvikda@buffalo.edu](mailto:souvikda@buffalo.edu))

TBD

### Course Details:

1. You are expected to attend all lectures and to complete all readings on time. Recordings will be made available shortly after live class concludes. The recordings are meant to serve as study aids, not as a substitute for attending class.
2. There will be 4 programming assignments in this course. The assignments cover the configuration of Solr for a particular search task, building of search indexes, evaluation of IR models, and a final (group) project requiring the development of a complete IR solution based on a real-world problem. All programming assignments will require the use of the Google Cloud Platform; more information on this will be provided in class.

3. We will use Piazza for course related discussion. The Piazza link is <https://piazza.com/buffalo/fall2022/cse435535>

Class notes will be posted there prior to class. Projects and announcements will also be posted on this site. Piazza should be used for Q&A related to the course and particularly projects.

**\*\*\* You should not post class materials (notes, exams, projects) on public sites: this would be a violation of Intellectual Property rights \*\*\*\***

4. Please read department policy on academic dishonesty; *this will be enforced strictly*.

UB Undergrad AI policy: [https://catalog.buffalo.edu/policies/academic\\_integrity\\_2019-20.html](https://catalog.buffalo.edu/policies/academic_integrity_2019-20.html)

UB Graduate AI policy: <https://grad.buffalo.edu/succeed/current-students/policy-library.academics.html#grievanceandintegrity>

CSE AI policy: <https://engineering.buffalo.edu/computer-science-engineering/information-for-students/policies/academic-integrity.html>

#### IMPORTANT DATES

First day of class	Aug 29
Midterm- 1	October 5
Midterm- 2	Nov 14
Final Project Presentation	Dec 7 & 9
Last Lecture	Dec 9
Project 1 Due	Sept 21
Project 2 Due	Oct 12
Project 3 Due	Nov 2
Project 4 Due	Dec 7

#### GRADING

Midterms	40%
Projects	60%
Total	100%

## COURSE SCHEDULE

Week and Date	Topics	Readings *	Key Activities
Week 1 Aug 29, 31, Sept 2	Introduction to IR Conceptual Models of IR Boolean Model <b>Project 1 release</b>	Chapter 1, 2	<ul style="list-style-type: none"> <li>Project 1 Release</li> <li>Create Reddit, GCP accounts</li> </ul>
Week 2 Sept 5 ( <b>holiday</b> ), 7, 9	Tokenization Text analysis: stop lists, stemming Dictionaries, Tolerant Retrieval	Chapter 3 Supplements	Recitation – SOLR, GCP setup (hands-on)
Week 3 Sept 12, 14, 16	Index Construction Distributed Indexing and Search Hadoop	Chapter 4 Supplements	
Week 4 Sept 19, 21, 23	Text Properties: Heaps, Zipfs Laws Index Compression Vector-Space Model <b>Project 2 release</b>	Chapter 5, 6	<ul style="list-style-type: none"> <li><b>Project 1 Due on Sept 21</b></li> <li>Project 2 Release</li> </ul>
Week 5 Sept 26, 28, 30	TF-IDF Weighting Scoring and Ranking in IR Systems	Chapter 6, 7	
Week 6 Oct 3, 5, 7	Evaluation Machine Learned Ranking <b>Midterm 1</b>	Chapter 8 Handouts	<b>Midterm 1</b>
Week 7 Oct 10, 12, 14	Relevance Feedback Query Expansion: Local and Global <b>Project 3 release</b>	Chapter 9	<ul style="list-style-type: none"> <li><b>Project 2 Due on Oct 12</b></li> <li>Project 3 Release</li> </ul>
Week 8 Oct 17, 19, 21	Probabilistic IR: Okapi (BM 25), DFR, Language Models	Chapter 11, 12	
Week 9 Oct 24, 26, 28	Prob IR contd. Text Classification	Chapter 13, 14	
Week 10 Oct 31, Nov 2, 4	Web Search Web Crawling	Chapter 19, 20	<ul style="list-style-type: none"> <li><b>Project – 3 Due on Nov 2</b></li> </ul>
Week 11 Nov 7, 9, 11	Social Network Analysis: Link Analysis, PageRank, HITS <b>Project 4 release</b>	Chapter 21 Handouts	<ul style="list-style-type: none"> <li>Project 4 Release</li> </ul>
Week 12 Nov 14, 16, 18	Word Vectors: Latent Semantic Indexing Word2Vec, GloVe, Doc2Vec  Using word embeddings in Search Computational Advertising	Chapter 18 Handouts	
Week 13 Nov 21	<b>Midterm 2</b>		<b>Midterm 2</b>
Nov 23, 25	<b>***THANKSGIVING BREAK***</b>		
Week 14 Nov 28, 30, Dec 2	E-commerce, social media search Knowledge Graphs	Handouts	
Week 15 Dec 5, 7, 9	<b>Student Project Presentations</b>		<b>Project – 4 Due on Dec 7</b>

\*Chapters are from the *An Introduction to Information Retrieval* textbook unless specified.