

[Information Retrieval for Data Science]

Course Administrative Details

Course Title	<i>Information Retrieval for Data Science</i>		
Instructor(s)	<i>Stanislav Protasov</i>	Instructor's e-mail	<i>s.protasov@innopolis.ru</i>
Course #	<i>XXX</i>	Course Type	<i>Core</i>
Faculty	<i>Computer Science and Engineering</i>	Major	<i>Data Science</i>
Academic year	<i>2019-2020</i>	Semester Offered	<i>Spring</i>
No. of Credits	<i>6 ECTS</i>	Total workload on average	<i>3 hours in auditory per week 3 hours homework per week 20 hours project</i>
Lecture Hours	<i>2 per week</i>	Lab Hours	<i>2 per week</i>
Language	<i>English</i>	Frequency	<i>Weekly</i>
Target Audience	<i>Masters</i>	Anticipated Enrollment	<i>30 students</i>
Studying year	<i>1 (5)</i>		
Grading Mode	<i>A, B, C, Fail</i>	Keywords	<i>Information retrieval, information search, search engines, data mining</i>

Course outline

The course is designed to prepare students to understand and learn contemporary tools of information retrieval systems. The course will focus on applying mathematical and programming tools to building such systems. Throughout the course, students will be involved in discussions, readings and assignments to experience real world systems. The technologies and algorithms covered in class includes advanced algorithms, linear algebra, machine learning, data mining, natural language processing and so on.

Course Delivery

The course will be given weekly from *January to April 2020*. Every week, there will be a 2-hour lecture followed by a 2-hour lab session. Lab sessions are followed by homework. In the end of the semester students at the exam present their projects – information retrieval services.

Prerequisite courses

Machine learning, Python, algorithms and data structures

Required background knowledge

Solid knowledge of imperative and object-oriented programming concepts and good programming skills in Python are assumed. Good understanding of machine learning together with a skill to run models required. Basic probability theory and statistics understanding is also required.

Course structure

[IA – Individual Assignment]

Week# / Date	Topic	Assignments
Week 1	<i>Introduction to information retrieval</i>	<i>Non-graded introduction test</i>
Week 2	<i>Building inverted index. Language, tokenization, stemming, searching, scoring.</i>	<i>Inverted index service with scoring</i>
Week 3	<i>Distributive semantics. Vector model. Dimension reduction.</i>	<i>Vector representation</i>

Week 4	<i>ML approaches to vector modelling.</i>	<i>ML for document representation</i>
Week 5	<i>Indexing for vector model. Kd-trees, quad trees, Annoy, FAISS, HNSW</i>	<i>Fast vector index. Recommender system</i>
Week 6	<i>Web basics. Internet crawling. XML and HTML processing. Dynamic documents processing.</i>	<i>Dynamic site parsing and robot ethics</i>
Week 7	<i>Spellchecking and query correction</i>	<i>Spellchecker</i>
Week 8	<i>Query expansion and suggest</i>	<i>Query expansion and suggest lab</i>
Week 9	<i>Language model. Topic model. Clustering and classification</i>	<i>Topic modeling</i>
Week 10	<i>On newsfeeds</i>	<i>Newsfeed</i>
Week 11	<i>Image and video processing. Understanding and enhancing</i>	<i>Image/Video to text</i>
Week 12	<i>Audio processing. Speech to text. Acoustic fingerprinting</i>	<i>Shazam</i>
Week 13	<i>Quality assessment. A/B testing, SBS, pFound, DGC, nDGC.</i>	<i>A/B test framework</i>
Week 14	<i>Web search specific topics. PageRank. Duplicates. CTR.</i>	<i>PageRank</i>
Week 15	<i>[Extension slot: selected topic]</i>	<i>--</i>

Textbook(s)

An Introduction to Information Retrieval by Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Cambridge University Press (any edition)

Reference Materials

TBD

Computer Resources

Students should have laptops with Anaconda installer.

Laboratory Exercises

Tutorial exercises will be set on weekly basis

Laboratory Resources

No laboratory resources are required for this course.

Grading criteria

Homework = 50%, Final Exam (project defense) = 50%

Late Submission Policy

No late submission allowed.

Cooperation Policy and Quotations

We encourage vigorous discussion and cooperation in this class. You should feel free to discuss any aspects of the class with any classmates. However, we insist that all assignments should be done by you alone. We will run automatic code comparison software on assignment submissions to compare solutions. Violations of this policy will be investigated and will result in zero scores on any assignments brought under suspicion.