



Information Retrieval Systems Software. Part 2.

Methods of Information Retrieval Organization

Syllabus

Requisites of the Course

Cycle of Higher Education	<i>First cycle of higher education (Bachelor's degree)</i>
Field of Study	<i>12 Information Technologies</i>
Speciality	<i>121 Software engineering</i>
Education Program	<i>Software Engineering of Multimedia and Information Retrieval Systems</i>
Type of Course	<i>Normative</i>
Mode of Studies	<i>full-time</i>
Year of studies, semester	<i>4 year (8 semester)</i>
ECTS workload	<i>4,5 credits (ECTS). Time allocation: 36 hours for lectures, 18 hours for lab classes, 81 hours for self-study</i>
Testing and assessment	<i>Exam, midterm test</i>
Course Schedule	<i>According to rozklad.kpi.ua</i>
Language of Instruction	<i>English</i>
Course Instructors	<i>Lecturer: teaching assistant, Yakiv Yusyn, yusyn@pzks.fpm.kpi.ua Teacher of laboratory work: teaching assistant, Yakiv Yusyn, yusyn@pzks.fpm.kpi.ua</i>
Access to the course	<i>Google Classroom. Link to be provided to registered students</i>

Outline of the Course

1. Course description, goals, objectives, and learning outcomes

The study of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" allows students to form the competencies necessary to solve practical problems of professional activities related to developing and using of information retrieval methods in software.

The **purpose** of studying the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" is to form in students the ability to develop software for information retrieval, as well as use third-party software and libraries for information retrieval.

The **subject** of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" is mathematical and algorithmic methods of presentation and search of unstructured documental information in information retrieval systems.

Study of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" contributes to the formation of **professional competencies (FC)** in students, necessary for solving practical tasks of professional activity, related to the development and use of information retrieval methods in software.

PC07 Knowledge of information data models, the ability to create software for data storage, retrieval and processing.

PC13 Ability to reasonably select and master software development and maintenance tools.

PC14 Ability to algorithmic and logical thinking.

PC15 Ability to apply fundamental and interdisciplinary knowledge to build advanced retrieval algorithms.

PC17 Ability to develop software for information retrieval systems.

PC20 Ability to apply the acquired fundamental mathematical knowledge to develop calculation methods in the multimedia and information retrieval systems creation.

Study of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" contributes to the formation in students of the following **program learning outcomes (PLO)** according to the educational program:

PLO03 To know the software life cycle basic processes, phases and iterations.

PLO05 To know and apply relevant mathematical concepts, domain methods, system and object-oriented analysis and mathematical modeling for software development.

PLO06 Ability to select and use the appropriate task of software development methodology.

PLO07 To know and to apply in practice the fundamental concepts, paradigms and basic principles of the functioning of language, instrumental and computational tools of software engineering.

PLO08 To know and to be able to develop a human-machine interface.

PLO10 To conduct a pre-project survey of the subject area, system analysis of the design object.

PLO11 To select initial data for design, guided by formal methods of describing requirements and modeling.

PLO12 To apply effective approaches to software design in practice.

PLO13 To know and apply methods of developing algorithms, designing software and data and knowledge structures.

PLO15 To choose programming languages and development technologies to solve the problems of creating and maintaining software.

PLO16 To have the software development, design approval and all types of software documentation release skills.

PLO17 To be able to apply methods of component software development.

PLO18 To know and be able to apply information technology of processing, storage and transmission of data.

PLO19 To know and be able to apply software verification and validation methods.

PLO20 To know approaches to evaluation and quality assurance of software.

PLO25 To know and to be able to use fundamental mathematical tools in the algorithms construction and modern software development.

PLO31 To be able to identify, analyze and document software requirements for multimedia and information retrieval systems.

PLO38 To be able to apply programming technologies for multimedia and information retrieval systems software development.

PLO39 To know the types of search engines, the principles of their construction, the methods and algorithms for performing different kinds of information retrieval in them.

PLO40 To know and be able to apply in practice the methods and criteria for estimating the effectiveness of information retrieval.

PLO42 To know the basic presentation models of textual and multimedia information and methods of its pre-processing for use in the design of information retrieval systems.

PLO43 To know and be able to use in practice the existing software resources and libraries for processing of textual information and multimedia data in information retrieval systems.

PLO44 To know the most common query languages used in the development of information retrieval systems.

2. Prerequisites and post-requisites of the course (the place of the course in the scheme of studies in accordance with curriculum)

Successful study of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" is preceded by the study of discipline "Information Retrieval Systems Software 1. NoSQL Databases" of the curriculum for Bachelors in 121 Software Engineering.

The theoretical knowledge and practical skills obtained during the mastering of the discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" ensure the successful implementation of course and diploma projects in the specialty 121 Software engineering. Also, the acquired knowledge and skills are a prerequisite for studying the disciplines "Information Retrieval Systems and Services" and "Artificial Intelligence Technologies for Information Retrieval Systems" of the Master's degree curriculum in 121 Software Engineering.

3. Content of the course

Discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" involves the study of the following topics:

Topic 1. Introduction to information retrieval

Topic 2. Presentation models of documents for information retrieval

Topic 3. Full-text search engine "Elasticsearch"

Topic 4. Overview of advanced tasks of information retrieval

Topic 5. Validation and verification of information retrieval results

Midterm test

Exam

4. Coursebooks and teaching resources

Main resources:

1. Manning C. D. *An Introduction to Information Retrieval* / C. D. Manning, P. Raghavan, H. Schütze. – Cambridge: Cambridge University Press, 2008. – p. 544.

Get acquainted with the sections related to the following topics of the discipline: introduction to information retrieval, presentation models of documents for information retrieval, overview of advanced tasks of information retrieval, validation and verification of information retrieval results. Open-Access Resource.

2. *Elasticsearch Guide* [Online]. – 2022. – Available: <https://www.elastic.co/guide/en/elasticsearch/reference/master/index.html>.

Get acquainted with the sections related to the following topics of the discipline: full-text search engine "Elasticsearch". Open-Access Resource.

3. *Kibana Guide* [Online]. – 2022. – Available: <https://www.elastic.co/guide/en/kibana/current/index.html>.

Get acquainted with the sections related to the following topics of the discipline: full-text search engine "Elasticsearch". Open-Access Resource.

Educational content

5. Methodology

No	Type of a class	A class description
----	-----------------	---------------------

<i>Topic 1. Introduction to information retrieval</i>		
<i>1</i>	<i>Lecture 1. Basic concepts of information retrieval</i>	<i>Introduction. The purpose and tasks of the discipline of information retrieval. The main conceptual apparatus of the discipline. Materials for self-studying: p. 6, No. 1.</i>
<i>2</i>	<i>Lecture 2. Information retrieval</i>	<i>The task of information retrieval. Problems of information retrieval. Main types of information retrieval. Applications of information retrieval. Requirements for information retrieval results. Materials for self-studying: p. 6, No. 2.</i>
<i>Topic 2. Presentation models of documents for information retrieval</i>		
<i>3</i>	<i>Lecture 3. Set-theoretic models</i>	<i>Standard Boolean model. Extended Boolean model. Fuzzy retrieval. Materials for self-studying: p. 6, No. 3.</i>
<i>4</i>	<i>Laboratory class 1. Implementation of presentation models (part 1)</i>	<i>Task: to implement a set-theoretic model. Materials for self-studying: p. 6, No. 4.</i>
<i>5</i>	<i>Lecture 4. Algebraic models. Part 1</i>	<i>Vector space model. TF-IDF. Materials for self-studying: p. 6, No. 5.</i>
<i>6</i>	<i>Lecture 5. Algebraic models. Part 2</i>	<i>Generalized vector space model. Topic-based vector space model. Latent semantic indexing. Latent semantic analysis. Materials for self-studying: p. 6, No. 6.</i>
<i>7</i>	<i>Laboratory class 2. Implementation of presentation models (part 2)</i>	<i>Task: to implement an algebraic model. Materials for self-studying: p. 6, No. 7.</i>
<i>8</i>	<i>Lecture 6. Probabilistic models. Part 1</i>	<i>Binary Independence Model. Probabilistic relevance model. Uncertain inference. Materials for self-studying: p. 6, No. 8.</i>
<i>9</i>	<i>Lecture 7. Probabilistic models. Part 2</i>	<i>Language models. Latent Dirichlet allocation. Materials for self-studying: p. 6, No. 9.</i>
<i>10</i>	<i>Laboratory class 3. Implementation of presentation models (part 3)</i>	<i>Task: to implement an algebraic model. Materials for self-studying: p. 6, No. 10.</i>
<i>Topic 3. Full-text search engine “Elasticsearch”</i>		
<i>11</i>	<i>Lecture 8. Full-text search engine “Elasticsearch”: overview, main concepts.</i>	<i>Elastic stack. Main properties of Elasticsearch engine. Glossary. Materials for self-studying: p. 6, No. 11.</i>
<i>12</i>	<i>Lecture 9. Indexing. Index creation. Filters and simple</i>	<i>Index creation in Elasticsearch, document storing in it. Difference between filters and queries. Filters</i>

	<i>queries using Elastic Query DSL.</i>	<i>execution. Simple queries. Materials for self-studying: p. 6, No. 12.</i>
13	<i>Laboratory class 4. Metadata search using Elasticsearch (part 1)</i>	<i>Task: to implement software for documents storing in Elasticsearch. Materials for self-studying: p. 6, No. 13.</i>
14	<i>Lecture 10. Complex queries using Elastic Query DSL</i>	<i>Complex queries: fuzzy, wildcard, regex, range. Materials for self-studying: p. 6, No. 14.</i>
15	<i>Lecture 11. Elasticsearch analyzers: stop-word removal, stemming</i>	<i>Analyzers in Elasticsearch. Stop-word definition. Stemming process. Standard Elasticsearch analyzer. Analyzer for English. Materials for self-studying: p. 6, No. 15.</i>
16	<i>Laboratory class 5. Metadata search using Elasticsearch (part 2)</i>	<i>Task: to add the possibility of metadata search using Elasticsearch filters. Materials for self-studying: p. 6, No. 16.</i>
17	<i>Lecture 12. Elasticsearch integration into existing software systems</i>	<i>Elasticsearch cluster. Best practices of Elasticsearch integration. Materials for self-studying: p. 6, No. 17.</i>
18	<i>Lecture 13. Data visualization tool Kibana. Query language KQL</i>	<i>Kibana as a tool for data visualization from Elasticsearch. Main concepts. Simple queries using Kibana Query Language. Materials for self-studying: p. 6, No. 18.</i>
19	<i>Laboratory class 6. Metadata search using Elasticsearch (part 3)</i>	<i>Task: to add the possibility of metadata search using Elasticsearch filters. Materials for self-studying: p. 6, No. 19.</i>
<i>Topic 4. Overview of advanced tasks of information retrieval</i>		
20	<i>Lecture 14. Documents classification</i>	<i>The task of documents classification. Methods of documents classification. Materials for self-studying: p. 6, No. 20.</i>
21	<i>Laboratory class 7. Full-text search using Elasticsearch (part 1)</i>	<i>Task: to implement software for full-text search using Elasticsearch. Materials for self-studying: p. 6, No. 21.</i>
22	<i>Lecture 15. Documents clustering</i>	<i>The task of documents clustering, the difference from classification. Methods of document clustering. Materials for self-studying: p. 6, No. 22.</i>
23	<i>Laboratory class 8. Full-text search using Elasticsearch (part 2)</i>	<i>Task: to add Elasticsearch analyzer. Materials for self-studying: p. 6, No. 23.</i>
<i>Topic 5. Validation and verification of information retrieval results</i>		
24	<i>Lecture 16. Criteria for evaluating the effectiveness of information retrieval</i>	<i>Criteria for evaluating the effectiveness: precision, recall, F-measure. Materials for self-studying: p. 6, No. 24.</i>
25	<i>Laboratory class 9. Full-text search using Elasticsearch (part 3)</i>	<i>Task: to add Elasticsearch analyzer. Materials for self-studying: p. 6, No. 25.</i>

26	<i>Lecture 17. Automated testing of information retrieval systems</i>	<i>Automated testing of information retrieval systems. Oracle problem. Metamorphic testing. Materials for self-studying: p. 6, No. 26.</i>
<i>Midterm test</i>		

6. Self-study

The discipline "Information Retrieval Systems Software 2. Methods of Organization of Information Retrieval" is based on independent preparation for classroom classes on theoretical and practical topics.

No	Topic name	Number of hours	Resources
1	<i>Preparation for the lecture 1</i>	1.5	1, preface
2	<i>Preparation for the lecture 2</i>	1.5	1, preface
3	<i>Preparation for the lecture 3</i>	1.5	1, pp. 1-16
4	<i>Preparation for the laboratory class 1</i>	2	All resources for lecture(s) No 3
5	<i>Preparation for the lecture 4</i>	1.5	1, pp. 109-132
6	<i>Preparation for the lecture 5</i>	1.5	1, pp. 109-132
7	<i>Preparation for the laboratory class 2</i>	2	All resources for lecture(s) No 4-5
8	<i>Preparation for the lecture 6</i>	1.5	1, pp. 219-251
9	<i>Preparation for the lecture 7</i>	1.5	1, pp. 219-251
10	<i>Preparation for the laboratory class 3</i>	2	All resources for lecture(s) No 4-5
11	<i>Preparation for the lecture 8</i>	1.5	2, sections "What is Elastic?", "Glossary"
12	<i>Preparation for the lecture 9</i>	1.5	2, section «Query DSL»
13	<i>Preparation for the laboratory class 4</i>	2	All resources for lecture(s) No 8-9
14	<i>Preparation for the lecture 10</i>	1.5	2, section "Query DSL"
15	<i>Preparation for the lecture 11</i>	1.5	2, section "Text analysis"
16	<i>Preparation for the laboratory class 5</i>	2	All resources for lecture(s) No 8-10
17	<i>Preparation for the lecture 12</i>	1.5	All resources for lecture(s) No 8-9
18	<i>Preparation for the lecture 13</i>	1.5	3, sections "What is Kibana?", "Quick Start", "Kibana Query Language"

19	<i>Preparation for the laboratory class 6</i>	2	<i>All resources for lecture(s) No 8-10, 12</i>
20	<i>Preparation for the lecture 14</i>	1.5	<i>1, pp. 289-307</i>
21	<i>Preparation for the laboratory class 7</i>	2	<i>All resources for lecture(s) No 8-12</i>
22	<i>Preparation for the lecture 15</i>	1.5	<i>1, pp. 349-395</i>
23	<i>Preparation for the laboratory class 8</i>	2	<i>All resources for lecture(s) No 8-12</i>
24	<i>Preparation for the lecture 16</i>	1.5	<i>1, pp. 151-172</i>
25	<i>Preparation for the laboratory class 9</i>	2	<i>All resources for lecture(s) No 8-12</i>
26	<i>Preparation for the lecture 17</i>	1.5	<i>All resources for lecture(s) No 12</i>
27	<i>Preparation for the midterm test</i>	7.5	<i>All resources for lecture(s) No 1-9, 14-16</i>
28	<i>Preparation for the exam</i>	30	<i>All resources for the semester</i>

Policy and Assessment

7. Course policy

- *Attendance at lectures is mandatory.*
- *Attendance at laboratory classes can be sporadic and if you need to defend a laboratory work.*
- *Rules of conduct in the classroom: activity, respect for those present.*
- *Adherence to the policy of academic integrity.*
- *Rules for protecting the laboratory work: the work should be done according to the option of the student, which is determined by his number in the group list.*
- *The rules for assigning penalty points are as follows.*

Penalty points are awarded for:

- *Plagiarism (completed task does not correspond to the variant of the task) in the laboratory work: - 5 points for each attempt.*

8. Monitoring and grading policy

During the semester, students complete 3 laboratory works. Maximum number of points for each work: 8 points

Points are awarded for:

- *quality of the work: 0-3 points;*
- *answer during the work defense of: 0-3 points;*
- *timely submission of work to the defense: 0-2 points.*

Criteria for evaluating the quality of work:

- 3 points – the work is done and performed qualitatively;*
- 2 points – the work is done qualitatively, but has shortcomings;*
- 1 point – the work is done qualitatively, but has some issues;*
- 0 points – the work is done qualitatively, but has serious issues.*

Answer evaluation criteria:

3 points – complete answer, well-reasoned;

2 points – complete answer, with minor issues;

1 point – complete answer, with major issues;

0 points – not complete or incorrect answer.

Timely submission of work to the defense:

2 points – the work is submitted for defense no later than the specified deadline;

0 points – the work is submitted for defense later than the specified deadline.

Maximum number of points for performing and defending laboratory works:

8 points × 3 laboratory works = 24 points.

The task for the midterm test consists of 18 test questions – 10 questions with one correct answer and 8 questions with few correct answers. Maximum number of points for each question with one correct answer is 1 point, maximum number of points for each question with few correct answers is 2 points.

Answer evaluation criteria for a question with one correct answer:

1 point – correct answer;

0 points – no answer or incorrect answer.

Answer evaluation criteria for a question with few correct answers:

2 points – all correct answers and no incorrect answers are selected;

1 point – at least 50% of all correct answers are selected;

0 points – no answer or all answers are wrong.

Maximum number of points for the midterm test:

1 point × 10 questions with one correct answer + 2 points × 8 questions with few correct answers = 26 points.

Semester component of the rating scale $R_C = 50$ points, it is defined as the sum of points received for the performance and defense of laboratory works and points received for the midterm test.

$R_C = 24 \text{ points} + 26 \text{ points} = 50 \text{ points}.$

Composition and evaluation criteria of the exam answer:

The exam consists of 3 questions – 2 theoretical and 1 practical.

Exam component of the rating scale $R_E = 50$ points.

Maximum number of points for the first theoretical question is 20 points, for the second theoretical question is 10 points, and for the practical question is 20 points.

Evaluation criteria for the first theoretical question:

18-20 points – complete answer, well-reasoned;

10-17 points – complete answer, with minor issues;

5-9 points – complete answer, with major issues;

0-4 points – not complete or incorrect answer.

Evaluation criteria for the second theoretical question:

9-10 points – complete answer, well-reasoned;

7-8 points – complete answer, with minor issues;

4-6 points – complete answer, with major issues;

0-3 points – not complete or incorrect answer.

Evaluation criteria for the practical question:

18-20 points – complete answer, well-reasoned;

10-17 points – complete answer, with minor issues;

5-9 points – complete answer, with major issues;

0-4 points – not complete or incorrect answer.

Maximum number of points for the exam:

$R_E = 20 + 10 + 20 = 50$ points.

The rating scale for the discipline is: $R = R_C + R_E = 50$ points + 50 points = 100 points

Semester assessment: Exam.

Prerequisite for admission to the exam is the performance and defense of all laboratory works.

The table of compliance between overall points and the final grade:

Points	Grade
100-95	Excellent
94-85	Very good
84-75	Good
74-65	Satisfactory
64-60	Fair
Less than 60	Unsatisfactory
Course requirements are not met	Not graded

Syllabus of the course

Is designed by teaching assistant, Yakiv Yusyn

Adopted by Computer Systems Software Department (protocol № 12 from 26.04.23)

Approved by the Faculty Board of Methodology (protocol № 10 from 26.05.23)