

## COURSE DESCRIPTION – ACADEMIC YEAR 2018/2019

<b>Course title</b>	<b>Data Curation</b>
<b>Course code</b>	73005
<b>Scientific sector</b>	ING-INF/05
<b>Degree</b>	Master in Computational Data Science (LM-18)
<b>Semester</b>	1
<b>Year</b>	1
<b>Credits</b>	12
<b>Modular</b>	Yes
<b>Total lecturing hours</b>	80
<b>Total lab hours</b>	40
<b>Attendance</b>	Attendance is not compulsory, but non-attending students have to contact the lecturers at the start of the course to agree on the modalities of the independent study.
<b>Prerequisites</b>	<p>Knowledge of relational databases, as taught in an introductory course at the BSc level. Basic knowledge of first-order logic, as taught in a BSc course in logic or discrete mathematics. Java programming skills for the project part.</p> <p>(Data profiling module): Basic knowledge of linear algebra, statistics and probability theory; MATLAB and Python basics for the lab session.</p>
<b>Course page</b>	<a href="https://ole.unibz.it/">https://ole.unibz.it/</a>

<b>Specific educational objectives</b>	<p>The course belongs to the type "caratterizzanti – discipline informatiche" in the curriculum "Data Analysis".</p> <p>The Data Integration module addresses a variety of problems related to the integration of heterogenous data sources, that range from structured data (such as relational databases), over semi-structured data (such as data on the Web, and tree- and graph-structured data), to unstructured (textual) data. It overviews foundational techniques for data integration, such as schema mappings, data and schema matching, and query processing in data integration, and does so considering different data representations that go beyond the relational model, such as RDF data, linked open data, and knowledge graphs. Architectures for data integration and data federation and their adoption to build comprehensive data integration solutions are studied. By attending the course, students will also learn how to design and build a data integration system, possibly exploiting existing data access and data federation technologies.</p> <p>The Data Profiling module considers the problem of profiling databases, in the general case where they refer to multimedia material (i.e., images, sounds, videos, etc). The course will cover the task of textual tag generation from audio-video material, to provide discriminative metadata. Standard techniques for detecting patterns and violations, dependencies, data normalization, similarity measures, and summarization will be carried out. Topics include: basics on image and sound analysis, multimedia discriminative dimensionality reduction, bag of X, textual tag extraction by supervised classification.</p>
----------------------------------------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

	Patterns and violations detection, dependencies detection, scrubbing and normalization. Similarity measures, duplicate detection, summary extraction.
--	-------------------------------------------------------------------------------------------------------------------------------------------------------

<b>Module 1</b>	<b>Data Integration</b>
<b>Module code</b>	73005A
<b>Module scientific sector</b>	ING-INF/05
<b>Lecturer</b>	<a href="#">Diego Calvanese</a>
<b>Contact</b>	POS 2.07, calvanese@inf.unibz.it, +39 0471 016160
<b>Scientific sector of lecturer</b>	ING-INF/05
<b>Teaching language</b>	English
<b>Office hours</b>	Announced on the webpage of the lecturer. During the lecture time span, in general Friday 16:00-18:00. Outside of the lecture time span, students are advised to confirm availability by email.
<b>Lecturing assistant (if any)</b>	--
<b>Contact LA</b>	--
<b>Office hours LA</b>	--
<b>Credits</b>	6
<b>Lecturing hours</b>	40
<b>Lab hours</b>	20
<b>List of topics</b>	<ul style="list-style-type: none"> <li>• Data integration architectures</li> <li>• Schema mapping</li> <li>• Data matching</li> <li>• Heterogeneous and web data</li> <li>• Data cleaning</li> <li>• Query processing for data integration</li> </ul>
<b>Teaching format</b>	Frontal lectures, exercises, and labs.

<b>Module 2</b>	<b>Data Profiling</b>
<b>Module code</b>	73005B
<b>Module scientific sector</b>	INF/01
<b>Lecturer</b>	Marco Cristani
<b>Contact</b>	Office n.104, marco.cristani@inf.unibz.it, +39 348 8619516
<b>Scientific sector of lecturer</b>	ING-INF/05
<b>Teaching language</b>	English
<b>Office hours</b>	Monday 10:00-11:30, Wednesday 13:00-14:00 (students are advised to arrange beforehand by email, in order to avoid waiting time).
<b>Lecturing assistant (if any)</b>	--
<b>Contact LA</b>	--
<b>Office hours LA</b>	--
<b>Credits</b>	6
<b>Lecturing hours</b>	40
<b>Lab hours</b>	20
<b>List of topics</b>	<ul style="list-style-type: none"> <li>• Detecting patterns and violations</li> <li>• Detecting dependencies</li> <li>• Scrubbing and normalization</li> <li>• Similarity measures</li> <li>• Duplicate detection</li> <li>• Summary extraction</li> </ul>
<b>Teaching format</b>	Frontal lectures and labs

<p><b>Learning outcomes</b></p>	<p>Knowledge and understanding:</p> <ul style="list-style-type: none"> <li>• D1.1 - Knowledge of the key concepts and technologies of data science disciplines</li> <li>• D1.2 - Understanding of the skills, tools and techniques required for an effective use of data science</li> <li>• D1.6 - Knowledge of the principles and methods of data curation</li> </ul> <p>Applying knowledge and understanding:</p> <ul style="list-style-type: none"> <li>• D2.1 - Practical application and evaluation of tools and techniques in the field of data science</li> <li>• 2.5 - Ability to apply, evaluate and develop methods and tools for the integration, cleaning, and quality of data</li> </ul> <p>Making judgments</p> <ul style="list-style-type: none"> <li>• D3.2 - Ability to autonomously select the documentation (in the form of books, web, magazines, etc.) needed to keep up to date in a given sector</li> </ul> <p>Communication skills</p> <ul style="list-style-type: none"> <li>• D4.1 - Ability to use English at an advanced level with particular reference to disciplinary terminology</li> <li>• D4.3 - Ability to structure and draft scientific and technical documentation</li> </ul> <p>Learning skills</p> <ul style="list-style-type: none"> <li>• D5.2 - Ability to autonomously keep oneself up to date with the developments of the most important areas of data science</li> </ul>
<p><b>Assessment</b></p>	<p>Oral exam and project work. The mark for each part of the exam is 18-30, or insufficient.</p> <p>The oral exam covers Module 1 "Data Integration" and Module 2 "Data Profiling", and comprises verification questions, and open questions to test knowledge application skills. It counts for 50% of the total mark.</p> <p>The project consists of two parts.        Part 1 covers Module 1 "Data Integration", and verifies whether the student is able to apply advanced data integration techniques and technologies taught or presented in the course to solve concrete problems. It is assessed through a final presentation, a demo, and a project report and can be carried out either individually or in a group of 2 students. It is discussed during the oral exam, and it counts for 25% of the total mark.        Part 2 covers Module 2 "Data Profiling", and consists in developing MATLAB/Python code, with the aim of checking how much the student has been able to implement and evaluate the tools and techniques discussed during the frontal lessons. It is discussed during the oral exam, and it counts for 25% of the total mark.</p>
<p><b>Assessment language</b></p>	<p>English</p>
<p><b>Assessment Typology</b></p>	<p>Monocratic</p>

<p><b>Evaluation criteria and criteria for awarding marks</b></p>	<p>The final mark is computed as the weighted average of the oral exam, Part 1 of the project, and Part 2 of the project. The exam is considered passed when all three marks are valid, i.e., in the range 18-30. Otherwise, the individual valid marks (if any) are kept for all 3 regular exam sessions, until also all others parts are completed with a valid mark. After the 3 regular exam sessions, all marks become invalid.</p> <p>Relevant for the oral exam: clarity of answers; ability to recall principles and methods, and deep understanding about the course topics presented in the lectures; skills in applying knowledge to solve exercises about the course topics; skills in critical thinking.</p> <p>Relevant for the project: skill in applying knowledge in a practical setting; ability to summarize in own words; ability to develop correct solutions for complex problems; ability to write a quality report; ability in presentation; ability to work in teams (for Part 1 of the project).</p> <p>Non-attending students have the same evaluation criteria and requirements for passing the exam as attending students.</p>
<p><b>Required readings</b></p>	<p>Required books:</p> <ul style="list-style-type: none"> <li>• A. Doan, A. Halevy &amp; Z. Ives (2012). Data Integration. Morgan Kaufmann. (ST270 D631)</li> <li>• Z. Abedjan, L. Golab, F. Naumann, &amp; T. Papenbrock (2018). Data Profiling. Synthesis Lectures on Data Management, 10(4), 1-154.</li> </ul> <p>Additional material (slides, notes of the lecturers) will be made available before each lesson.</p> <p>Subject Librarian: David Gebhardi, <a href="mailto:David.Gebhardi@unibz.it">David.Gebhardi@unibz.it</a></p>
<p><b>Supplementary readings</b></p>	<p>Suggested readings (i.e., only a few chapters will be covered)</p> <ul style="list-style-type: none"> <li>• R. C. Gonzalez, &amp; R. E. Woods (2007). Image processing. Digital image processing, 2, 1.</li> <li>• R. O. Duda, P. E. Hart &amp; D. G. Stork (2012). Pattern classification. John Wiley &amp; Sons.</li> </ul>
<p><b>Software used</b></p>	<ul style="list-style-type: none"> <li>• Ontop system for ontology-based data access developed by the In2Data research group at the Faculty of Computer Science</li> <li>• Data federation tools</li> <li>• MATLAB 2017 (Image processing toolbox, statistic toolbox)</li> <li>• Python3.5</li> <li>• Pycharm 2016.3.2</li> </ul>