

Syllabus

Kursbeschreibung

Titel der Lehrveranstaltung	Data Curation
Code der Lehrveranstaltung	73056
Zusätzlicher Titel der	
Lehrveranstaltung	
Wissenschaftlich-	
disziplinärer Bereich	
Sprache	Englisch
Studiengang	Master in Computing for Data Science
Andere Studiengänge (gem. Lehrveranstaltung)	
Dozenten/Dozentinnen	Prof. Diego Calvanese,
	Diego.Calvanese@unibz.it
	https://www.unibz.it/en/faculties/engineering/academic-
	staff/person/3562
	Dr. Alessandro Mosca,
	Alessandro.Mosca@unibz.it
	https://www.unibz.it/en/faculties/engineering/academic-
	staff/person/29044
Wissensch.	
Mitarbeiter/Mitarbeiterin	
Semester	Erstes Semester
Studienjahr/e	2
KP	12
Vorlesungsstunden	80
Laboratoriumsstunden	40
Stunden für individuelles	180
Studium	
Vorgesehene Sprechzeiten	
Inhaltsangabe	Data integration architectures
	Query processing in data integration
	Schema mapping

- Data integration via virtual knowledge graphs
- Schema matching
- Data and entity matching
- Patterns and violation detection
- Correlation coefficients
- Elementary data analysis
- Exact and approximate value dependencies detection
- Normalization via data profiling
- Non-relational data profiling
- Data profiling for query optimisation

Themen der Lehrveranstaltung

The Data Preparation and Integration module addresses a variety of problems related to the integration of heterogenous data sources. It overviews the main issues in data integration, notably handling different forms of heterogeneity, and presents the general architecture of data integration systems. Foundational techniques for data integration are covered, such as data matching, schema matching and mapping, and query processing in data integration. A specific data integration approach relying on the technology of Virtual Knowledge Graphs and semantic mappings is presented in detail. The integration both of relational data sources, and of other types of data sources accessed by relying on data federation technology are considered. By attending the course, students will learn how to design and build a comprehensive data integration solution, possibly exploiting existing data access and data federation technologies.

The Data Profiling module considers a variety of problems related to the profiling of relational and non-relational data. It first overviews techniques for the exploratory analysis of relational databases, which are also adopted in preparatory activities of data cleansing and integration, and mining. Widely used statistical measures, such as the Pearson and Spearman correlation coefficients, are discussed and applied to real word examples in the first part of the course. Special attention will be also paid to the Phi-K correlation coefficient, which has been recently proposed to overcome well known limitations of the former two, such as the possibility to analyse correlations between numerical, categorical, interval and ordinal variables. Algorithms for discovering patterns and relational dependencies in tabular data will be presented in detail. In particular, the main part of the course will go through the

rmal specifications of uniqueness, functional and inclusion pendencies in relational data, with a special focus on state-of-e-art algorithms and their complexity. If time permits, the last rt of the course will discuss the application of the introduced of the course will discuss the application of the introduced of the course will discuss the application of the introduced of the course will discuss the application of the introduced of the course will discuss the application problem, do to a high level introduction of recent profiling methods that in be applied to RDF data. In a Integration, Schema Mappings, , Virtual Knowledge Graphs, elational Dependencies In the Data Preparation and Integration module: Knowledge of lational databases, as taught in an introductory course at the BSC rel. Basic knowledge of first-order logic, as taught in a BSC rurse in logic or discrete mathematics. Knowledge of Java or thon for the project part. In the Data Profiling module: relational database concepts, basic archine learning concepts; good to have knowledge of basic formation retrieval concepts. Python programming basics for the
r the Data Preparation and Integration module: Knowledge of lational databases, as taught in an introductory course at the BSc vel. Basic knowledge of first-order logic, as taught in a BSc urse in logic or discrete mathematics. Knowledge of Java or thon for the project part. In the Data Profiling module: relational database concepts, basic achine learning concepts; good to have knowledge of basic formation retrieval concepts. Python programming basics for the
lational databases, as taught in an introductory course at the BSc vel. Basic knowledge of first-order logic, as taught in a BSc urse in logic or discrete mathematics. Knowledge of Java or thon for the project part. In the Data Profiling module: relational database concepts, basic achine learning concepts; good to have knowledge of basic formation retrieval concepts. Python programming basics for the
o sessions.
ontal lectures, exercises, and labs.
tendance is not compulsory, but non-attending students must ntact the lecturers at the start of the course to agree on the odalities of the independent study.
The course belongs to the type "caratterizzanti – discipline formatiche" in the curriculum "Data Analysis". The Data Integration module addresses a variety of problems lated to the integration of heterogenous data sources. It rerviews the main issues in data integration, notably handling referent forms of heterogeneity, and presents the general chitecture of data integration systems. Foundational techniques of data integration are covered, such as data matching, schema atching and mapping, and query processing in data integration. A recific data integration approach relying on the technology of trual Knowledge Graphs and semantic mappings is presented in

learn how to design and build a comprehensive data integration solution, possibly exploiting existing data access and data federation technologies.

The Data Profiling module considers a variety of problems related to the profiling of relational and non-relational data. It first overviews techniques for the exploratory analysis of relational databases, which are also adopted in preparatory activities of data cleansing and integration, and mining. Widely used statistical measures, such as the Pearson and Spearman correlation coefficients, are discussed and applied to real word examples in the first part of the course. Special attention will be also paid to the Phi-K correlation coefficient, which has been recently proposed to overcome well known limitations of the former two, such as the possibility to analyse correlations between numerical, categorical, interval and ordinal variables. Algorithms for discovering patterns and relational dependencies in tabular data will be presented in detail. In particular, the main part of the course will go through the formal specifications of uniqueness, functional and inclusion dependencies in relational data, with a special focus on state-ofthe-art algorithms and their complexity. If time permits, the last part of the course will discuss the application of the introduced profiling techniques to the relational query optimization problem, and to high level introduction of recent profiling methods that can be applied to RDF data.

Knowledge and understanding:

- D1.1 Knowledge of the key concepts and technologies of data science disciplines
- D1.2 Understanding of the skills, tools and techniques required for an effective use of data science
- D1.6 Knowledge of the principles and methods of data curation

Applying knowledge and understanding:

- D2.1 Practical application and evaluation of tools and techniques in the field of data science
- 2.5 Ability to apply, evaluate and develop methods and tools for the integration, cleaning, and quality of data
- 2.10 Application of languages, tools, and methods for the design of information systems and their corresponding software applications for data, process, and organization management



	 Making judgments 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 Communication skills D4.1 - Ability to use English at an advanced level with particular reference to disciplinary terminology D4.3 - Ability to structure and draft scientific and technical documentation Learning skills D5.2 - Ability to autonomously keep oneself up to date with the developments of the most important areas of data science
Spezifisches Bildungsziel und erwartete Lernergebnisse (zusätzliche Informationen)	
Art der Prüfung	Oral exam and project work. The mark for each part of the exam is 18-30, or insufficient. 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.
	The project consists of two parts. Part 1 covers Module 1 "Data Preparation and Integration" and verifies whether the student is able to apply advanced data integration techniques and technologies taught or presented in the course to solve a concrete integration problem. 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 verifies whether the student is able to apply data profiling techniques to a concrete use case. It is assessed through a final presentation 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.

Bewertungskriterien	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 other parts are completed with a valid mark. After the 3 regular exam sessions, all marks become invalid.
	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.
	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.
	Non-attending students have the same evaluation criteria and requirements for passing the exam as attending students.
Pflichtliteratur	A. Doan, A. Halevy & Z. Ives (2012). Data Integration. Morgan Kaufmann. (ST270 D631)
	Z. Abedjan, L. Golab, F. Naumann, & T. Papenbrock (2018). Data Profiling. Synthesis Lectures on Data Management, 10(4), 1-154.
	Z. Zhang, R. Zhang (2008). Multimedia Data Mining: A Systematic Introduction to Concepts and Theory. CRC Press LLC.
	Additional material (slides, notes of the lecturers) will be made available before each lesson.
	Subject Librarian: David Gebhardi, <u>David.Gebhardi@unibz.it</u>
Weiterführende Literatur	R. C. Gonzalez, & R. E. Woods (2007). Image processing. Digital image processing, 2, 1.

	R. O. Duda, P. E. Hart & D. G. Stork (2012). Pattern classification. John Wiley & Sons.
Weitere Informationen	Software used:
	Ontop system for deploying and querying Virtual Knowledge Graphs, developed by the In2Data research group at the Faculty of Computer Science.
	Relational DBMS, such as PostgreSQL.
	Data federation tools such as Denodo, Dremio, Teiid.
	Python3.5 with pandas, scikit-learn, scikit-image.
Ziele für nachhaltige Entwicklung (SDGs)	Hochwertige Bildung

Kursmodul

Titel des Bestandteils der	Data Preparation and Integration
Lehrveranstaltung	
Code der Lehrveranstaltung	73056A
Wissenschaftlich-	ING-INF/05
disziplinärer Bereich	
Sprache	Englisch
Dozenten/Dozentinnen	Prof. Diego Calvanese,
	Diego.Calvanese@unibz.it
	https://www.unibz.it/en/faculties/engineering/academic-
	staff/person/3562
Wissensch.	
Mitarbeiter/Mitarbeiterin	
Semester	Erstes Semester
KP	6
Verantwortliche/r Dozent/in	
Vorlesungsstunden	40
Laboratoriumsstunden	20

Stunden für individuelles Studium	90
Vorgesehene Sprechzeiten	
Inhaltsangabe	 Data integration architectures Query processing in data integration Schema mapping Data integration via virtual knowledge graphs Schema matching Data and entity matching
Themen der Lehrveranstaltung	The Data Preparation and Integration module addresses a variety of problems related to the integration of heterogenous data sources. It overviews the main issues in data integration, notably handling different forms of heterogeneity, and presents the general architecture of data integration systems. Foundational techniques for data integration are covered, such as data matching, schema matching and mapping, and query processing in data integration. A specific data integration approach relying on the technology of Virtual Knowledge Graphs and semantic mappings is presented in detail. The integration both of relational data sources, and of other types of data sources accessed by relying on data federation technology are considered. By attending the course, students will learn how to design and build a comprehensive data integration solution, possibly exploiting existing data access and data federation technologies.
Unterrichtsform	Frontal lectures, exercises, and labs.
Pflichtliteratur	A. Doan, A. Halevy & Z. Ives (2012). Data Integration. Morgan Kaufmann. (ST270 D631) Z. Abedjan, L. Golab, F. Naumann, & T. Papenbrock (2018). Data Profiling. Synthesis Lectures on Data Management, 10(4), 1-154. Z. Zhang, R. Zhang (2008). Multimedia Data Mining: A Systematic Introduction to Concepts and Theory. CRC Press LLC. Additional material (slides, notes of the lecturers) will be made available before each lesson.



Weiterführende Literatur	R. C. Gonzalez, & R. E. Woods (2007). Image processing. Digital image processing, 2, 1.	
	R. O. Duda, P. E. Hart & D. G. Stork (2012). Pattern classification. John Wiley & Sons.	

Kursmodul

Titel des Bestandteils der	Data Profiling
Lehrveranstaltung	
Code der Lehrveranstaltung	73056B
Wissenschaftlich-	INF/01
disziplinärer Bereich	
Sprache	Englisch
Dozenten/Dozentinnen	Dr. Alessandro Mosca,
	Alessandro.Mosca@unibz.it
	https://www.unibz.it/en/faculties/engineering/academic-
	staff/person/29044
Wissensch.	
Mitarbeiter/Mitarbeiterin	
Semester	Erstes Semester
KP	6
Verantwortliche/r Dozent/in	
Vorlesungsstunden	40
Laboratoriumsstunden	20
Stunden für individuelles	90
Studium	
Vorgesehene Sprechzeiten	
Inhaltsangabe	Patterns and violation detection
	Correlation coefficients
	Elementary data analysis
	Exact and approximate value dependencies detection
	Normalization via data profiling
	Non-relational data profiling
	Data profiling for query optimisation

Themen der	The Data Profiling module considers a variety of problems related
Lehrveranstaltung	to the profiling of relational and non-relational data. It first overviews techniques for the exploratory analysis of relational databases, which are also adopted in preparatory activities of data cleansing and integration, and mining. Widely used statistical measures, such as the Pearson and Spearman correlation coefficients, are discussed and applied to real word examples in the first part of the course. Special attention will be also paid to the Phi-K correlation coefficient, which has been recently proposed to overcome well known limitations of the former two, such as the possibility to analyse correlations between numerical, categorical, interval and ordinal variables. Algorithms for discovering patterns and relational dependencies in tabular data will be presented in detail. In particular, the main part of the course will go through the formal specifications of uniqueness, functional and inclusion dependencies in relational data, with a special focus on state-of-the-art algorithms and their complexity. If time permits, the last part of the course will discuss the application of the introduced profiling techniques to the relational query optimization problem, and to high level introduction of recent profiling methods that can be applied to RDF data.
Unterrichtsform	Frontal lectures and labs.
Pflichtliteratur	 A. Doan, A. Halevy & Z. Ives (2012). Data Integration. Morgan Kaufmann. (ST270 D631) Z. Abedjan, L. Golab, F. Naumann, & T. Papenbrock (2018). Data Profiling. Synthesis Lectures on Data Management, 10(4), 1-154. Z. Zhang, R. Zhang (2008). Multimedia Data Mining: A Systematic Introduction to Concepts and Theory. CRC Press LLC. Additional material (slides, notes of the lecturers) will be made available before each lesson.
Weiterführende Literatur	R. C. Gonzalez, & R. E. Woods (2007). Image processing. Digital image processing, 2, 1.



John Wiley & Sons.
