

# Advanced Data Science elective package

Advance Data Science	
<b>Offered by</b>	Department of Mathematics and Computer Science
<b>Language</b>	English
<b>Primarily interesting for</b>	All students, but most relevant for students with background in BSc in Data Science (DS)
<b>Prerequisites</b>	Required courses: Students are assumed to have basic skills in logic, set theory, calculus, discrete mathematics, databases, linear algebra, programming, algorithms, and data mining and machine learning. Recommended courses: -
<b>Contact person</b>	dr. Dirk Fahland

## Content and composition

Analysis of information, data, and knowledge is increasingly important, with broad application across science, engineering, society, and industry. To tackle these challenges, a wide variety of knowledge and skills in managing, mining, and analyzing of (big) data collections is necessary. This elective package provides deeper study of the foundations and applications of analysis of data as well as the data-driven construction and analysis of various kinds of models.

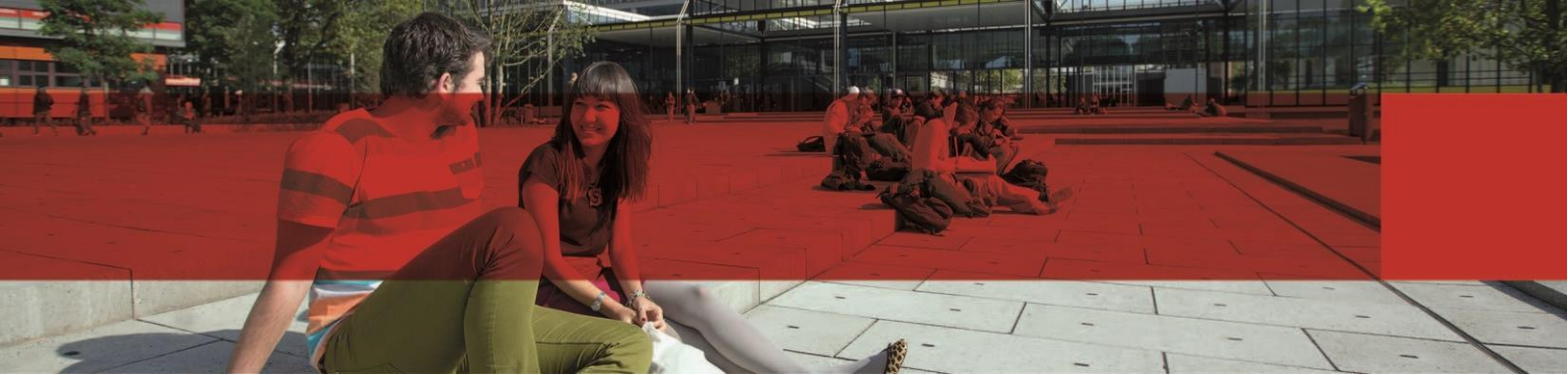
Students should select 3 out of the following package of courses. The courses can be followed in any order, but need to take the specific prerequisites for each course into account when scheduling this package. For example, 2ID70 has as prerequisite JBI050 in Q2/B.

Course code	Course name	Level classification
JBI060	Fundamentals of Process Mining	2.
2ID70	Data-intensive systems and applications	3.
2IX30	Responsible data science	3.
JBM170	Field Data Acquisition and Analysis	2.
JBM035	Linear Optimization for Data Science	2.
JBM090	Combinatorial Optimization for Data Science	3.

## Course description

### Fundamentals of Process Mining (JBI060)

In organizations, work is typically organized in the form of processes. Requests, orders, offers, invoices, or in general 'cases', are all handled by organizations for their customers or clients according to pre-defined activities executed by their staff. Event data records multi-variate event sequences of who executed which activity at which point in time for which case along with additional attributes relevant to the process case. The sequence of events recorded for a case reflects the decisions made on which activities need to be performed. Such event data is omnipresent and can be used to gain insights into the efficiency of an organization as well as insights into the compliance of processes to quality requirements and regulations. The use of traditional machine learning or data mining techniques on such sequential event-data is not trivial and this is the focus of the research area of process mining.



In this course, students will get acquainted with fundamental data science concepts and tasks in the context of processes. Students will learn how to encode sequences of events, how to use them for classical data science tasks such as prediction and recommendation and students will get introduced to process mining, a set of tools and techniques to visualize process data. Process mining uses process models another kind of models that are different from typical data mining and machine learning models. In the course, students will learn the relation between event data and activities, process models, patterns, and decision points. The course uses a set of process mining techniques implemented in the Python framework PM4Py, combined with an industrial process mining tool (through an academic license).

### **Data-intensive systems and applications (2ID70)**

This course prepares students to meet the new challenges of contemporary data engineering in which traditional assumptions break, where new data models, query languages and programming interfaces are required. In this course, we study how traditional relational database techniques such as indexing, query planning and optimization, transaction management and self-tuning can be made to work on a massive scale of thousands of machines and petabytes of data. We study models of contemporary data-intensive systems, their efficient engineering, and their practical use. These models include scalable data processing platforms (e.g., MapReduce, Spark) and stream processing engines. We discuss why these models were introduced, their relative advantages and disadvantages, how they are engineered, and how to effectively use them in practice.

### **Responsible Data Science (2IX30)**

The course is focused on studying the problems of fairness, accountability, confidentiality, and transparency (FACT) in data science, and data mining and machine learning in particular. One important challenge to face is that machine learnt models typically are not 100% accurate, i.e. in some ways these models are wrong. Thus, it is important to study how we can make a good use of models that are not perfect, how we can understand the strengths and weaknesses of these models, how we can help a decision maker to trust (or not trust) the model or its particular prediction, and how we can get insights into impact of input features and some inner logic of a predictive model. We need techniques not just to explain the decision of a model, but also to uncover and characterize undesired or even unlawful biases in its performance. Hence, the other important challenge to study is how to formally define such biases, how to uncover and quantify them and how to design machine learning solutions that would enable the so-called fair algorithmic decision making by design. On the other side of the spectrum, there are challenges of privacy and confidentiality. We will study the main principles and techniques that have been researched and employed in data mining for privacy preserving and secure computation to induce models from data and to apply them in real-life scenarios.

### **Field Data Acquisition and Analysis (JBM170)**

In this course students will learn: acquiring, analyzing and representing data in a real-life setting on energy intake (food/nutrition) and energy expenditure (active lifestyle)

How to handle the data acquired by 24/7 in a real-life field context?

Who has access to this data?

How to present/represent it in a meaningful manner?

The focus will be partly in the biophysical domain (food/activities), the technical domain (data acquisition and analytics) and partly in the domain of (interaction) psychology (behavioral change, contextual and adaptive data representation). Next to this also the corresponding legal and ethical aspects will be discussed.

### **Linear Optimization for Data Science (JBM035)**

Linear optimization is one of the fundamental computational tools in Operations Research, and is used for airline scheduling, production planning, and in many other industrial settings. In fact, it has been called one of the mathematical problems "using up most of the computer time in the world". We will first look at how linear optimization models arise from practical decision problems. Next we will consider the links with linear algebra and geometry. This will lead us to an algorithm, called the simplex method, that may be implemented using techniques from linear algebra. Every linear optimization problem has an associated dual problem, that has an economic interpretation in terms of "shadow prices", and we will look at these ideas in some detail. Finally, we will consider the case where the decision variables should take integer values, and study a techniques for such problems, namely branch-and-bound.



# Advanced Data Science elective package

## **Combinatorial Optimization for Data Science (JBM090)**

Combinatorial optimization consists in finding an optimal solution among a finite set of possible solutions. The main objectives of the course are to understand the complexity of different combinatorial optimization problems, model them with Integer Linear Programming (ILP), and learn different solution approaches to solve them. For some problems (the 'easy' ones) fast solution methods exist while for the other problems (the 'difficult' ones) finding the optimal solution is time consuming and can only be achieved in practice if the problem is 'small'. For the difficult problems one may have to use 'heuristics', which apply "rules of thumb" to find feasible solutions. Heuristics usually lead to 'good' quality solutions but without any guarantee on their optimality. In this course we will treat the aspects mentioned above as well as a number of concrete problems and solution methods. Among others, we will treat the easy problems 'shortest path', 'maximal flow' and 'assignment', the difficult problems 'traveling salesman', 'knapsack', and 'vehicle routing', and the methods 'greedy', 'branch and bound', 'dynamic programming', and 'local search'. We will also see how to implement some of these methods in Python.