

# Knowledge Discovery in Datenbanken I (IN5042)

---

Titel	Knowledge Discovery in Databases I	
Typ	Vorlesung mit Übung	
Credits	6 ECTS	
Lehrform/SWS	3V + 2Ü	
Sprache	Deutsch	
Modulniveau	Master	
Arbeitsaufwand	Präsenzstunden	75 Stunden
	Eigenstudium	105 Stunden
	Gesamtaufwand	180 Stunden
Angestrebte Lernergebnisse	<p>Die Studierenden kennen den grundlegenden Prozess des Knowledge Discovery in Datenbanken und die einzelnen Schritte dieses Prozesses sowie die grundlegenden Problemstellung im Data Mining (wie z.B. Feature Transformation, Suche nach häufigen Mustern, Musterevaluation). Sie sind in der Lage, Merkmalsräume, Ähnlichkeitsmaße und Distanz-Metriken zu beschreiben, zu analysieren, zu bewerten und gezielt anzuwenden. Sie sind in der Lage, grundlegende Verfahren (wie z.B. Clustering-Methoden, Klassifikatoren, Outlier Detection Methoden, Regressionsverfahren) in verschiedene Bereiche des Data Mining zu beurteilen, auszuwählen und gezielt einzusetzen sowie die gefundenen Muster und Funktionen zu evaluieren und kritisch zu interpretieren. Sie sind in der Lage, für ein gegebenes Problem (wie z.B. Erkennen von Spam Emails, Warenkorbanalyse, Clustern von Bildinhalten) einen Knowledge Discovery Prozess zu entwerfen und umzusetzen sowie aus den vorhanden Verfahren geeignete Verfahren auszuwählen und an diese Probleme anzupassen.</p>	
Intended Learning Outcomes	<p>Students are able to reproduce the process of knowledge discovery and fundamental problems of data mining (e.g. feature transformation, frequent pattern mining, pattern evaluation). They are able to describe, to analyze, to evaluate and to apply feature spaces, similarity measures, and distance metrics. They are able to employ and implement methods (e.g. clustering methods, classification methods, regression methods, outlier detection methods) for data mining tasks and to</p>	

	<p>evaluate the computed patterns and functions. They are able to design and implement knowledge discovery processes for given problem settings (e.g. spam detection, market basket analysis, or image clustering) and to select the best suited data mining methods for these problem settings.</p>
<p>Inhalt</p>	<p>Knowledge Discovery and Data Mining:</p> <ul style="list-style-type: none"> <li>• Definition Knowledge Discovery und Data Mining</li> <li>• Der KDD Prozess (einzelne Schritte, iterativer Ablauf)</li> <li>• Supervised und Unsupervised Learning</li> <li>• Grundliegende Aufgaben des Data Mining: Klassifikation, Clustering, Outlier Detection, Regression, Frequent Pattern Mining.</li> </ul> <p>Merkmalsräume:</p> <ul style="list-style-type: none"> <li>• Wahrscheinlichkeitsverteilungen (einfache univariate und multivariate Verteilungen, Abhängigkeit von Zufallsvariablen)</li> <li>• Distanzmaße und Ähnlichkeitsmaße (mathematische Eigenschaften wie Reflexivität, Symmetrie, Transitivität)</li> <li>• Beispiele für Feature-Transformationen (z.B. Farbhistogramme, Bag of Words)</li> <li>• einfache Verfahren zu Feature Selection (z.B. greedy forward selection) [optional]</li> <li>• einfache Verfahren der Feature Reduction (z.B. PCA) [optional]</li> </ul> <p>Klassifikation:</p> <ul style="list-style-type: none"> <li>• Evaluation von Klassifikatoren (Testschemata z.B. Crossvalidation, Bootstrapping, leaveone-out, Metriken )</li> <li>• Formale Aspekte des Lernens (Generalisierung, Overfitting, Problemdefinition)</li> <li>• Entscheidungsbäume</li> <li>• Bayes-Klassifikationen (naive Bayes, Bayes-Netze, diskrete und kontinuierliche Verteilungen)</li> <li>• Instanzbasierte Klassifikation.</li> <li>• fortgeschrittene Klassifikationsverfahren ( z.B. Support Vector Maschinen, Neuronale Netze, Gauss Klassifikatoren, logistische Regression) [optional]</li> <li>• regelbasierte Klassifikation und Inductive logical programming [optional]</li> <li>• Deep Learning Methoden [optional]</li> </ul> <p>Regression:</p> <ul style="list-style-type: none"> <li>• Problemdefinition (Bewertung von</li> </ul>

	<p>Regressionsmodellen)</p> <ul style="list-style-type: none"> <li>• einfache lineare Regressionsmodelle</li> <li>• Grundlegende Verfahren der multivariaten Regression</li> <li>• fortgeschrittene Regressionsverfahren (z.B. kernelbasierte Regression, instanzbasierte Regression).</li> </ul> <p>Clustering:</p> <ul style="list-style-type: none"> <li>• Problemdefinition (Zielsetzung, Abgrenzung zur Klassifikation)</li> <li>• Partitionierende Clusteringmethoden (k-Means, Expectation Maximization, weitere Verfahren z.B. PAM, CLARANCE, k-Modes)</li> <li>• Dichtebasiertes und hierarchische Clustering (z.B. DBSCAN, OPTICS, Single Link)</li> <li>• Self Organizing Maps [optional]</li> <li>• graphbasiertes Clustering und Spectral Clustering [optional]</li> <li>• Clusterevaluation [optional]</li> </ul> <p>Outlier Detection:</p> <ul style="list-style-type: none"> <li>• Aufgabenstellung (verschiedene Outlier Definitionen, Abgrenzung zu Clustering und Klassifikation)</li> <li>• statistische Outlier</li> <li>• distanzbasierte Outlierfaktoren</li> <li>• lokale Outlier (z.B. LOF).</li> <li>• fortgeschrittene Verfahren (z.B. ABOD) [optional]</li> <li>• Evaluation von Outlierverfahren [optional]</li> </ul> <p>Frequent Itemset Mining und Assoziationsregeln:</p> <ul style="list-style-type: none"> <li>• Einführung Pattern Mining (Häufigkeit, Konfidenz, Monotonie)</li> <li>• Frequent Itemset Mining (Suchraum, Apriori)</li> <li>• Assoziationsregeln (Ableitung, Interessantheit)</li> <li>• weiterführende Algorithmen zur Berechnung von frequent Itemsets [optional]</li> <li>• Datenstrukturen zur Suche in frequent Itemsets [optional]</li> </ul>
Contents	<p>Knowledge Discovery and Data Mining:</p> <ul style="list-style-type: none"> <li>• Definition Knowledge Discovery and Data Mining</li> <li>• KDD Process (different steps, iterative approach)</li> <li>• Supervised and unsupervised learning</li> <li>• Basic Data Mining tasks: Classification, Clustering, Outlier Detection, Regression, Frequent Pattern Mining.</li> </ul> <p>Feature Spaces:</p>

	<ul style="list-style-type: none"> <li>• Probability distributions (simple univariate and multivariate distributions, dependency of random variables)</li> <li>• Distance and similarity measures (mathematical characteristics such as reflexivity, symmetry, transitivity)</li> <li>• Examples for simple feature transformations (e.g. color histograms, bag of words)</li> <li>• simple methods for feature selection (e.g. greedy forward selection) [optional]</li> <li>• simple methods for feature reduction (e.g. PCA) [optional]</li> </ul> <p>Classification:</p> <ul style="list-style-type: none"> <li>• Classifier evaluation (testing schemes e.g. cross validation, bootstrapping, leave-one-out, evaluation metrics)</li> <li>• Formal aspects of learning (generalization, overfitting)</li> <li>• Decision trees</li> <li>• Bayes classifier (naive Bayes, Bayesian networks)</li> <li>• instance based Classification</li> <li>• advanced classification methods (e.g. support vector machines, neuronal Networks, Gaussian classifiers, logistic regression) [optional]</li> <li>• rule-based classifiers and inductive logical programming [optional]</li> <li>• deep learning [optional]</li> </ul> <p>Regression:</p> <ul style="list-style-type: none"> <li>• Problem definition (Evaluation of regression functions)</li> <li>• Simple linear regression</li> <li>• Basic methods for multivariate regression</li> <li>• Advanced regression methods (e.g. kernel based regression, instance-based regression)</li> </ul> <p>Clustering:</p> <ul style="list-style-type: none"> <li>• Problem definition (aims, difference to classification)</li> <li>• Partitioning clustering methods ( k-Means, expectation maximization, further methods e.g. PAM, CLARANCE, k-Modes)</li> <li>• Density-based and hierarchical clustering (e.g. DBSCAN, OPTICS, Single Link)</li> <li>• Self organizing maps [optional]</li> <li>• Graph-based clustering and spectral clustering [optional]</li> <li>• Evaluation of clusterings</li> </ul> <p>Outlier Detection:</p>
--	--

	<ul style="list-style-type: none"> <li>• General setting (various outlier definitions, differences to clustering and classification)</li> <li>• statistic outliers</li> <li>• distance-based outliers</li> <li>• local outlier (e.g. LOF)</li> <li>• Advanced methods for outlier detection (e.g. ABOD) [optional]</li> <li>• Evaluation of outlier detection methods [optional]</li> </ul> <p>Frequent Itemset Mining and Association Rules:</p> <ul style="list-style-type: none"> <li>• Introduction to Pattern Mining (Frequency, Confidence, Monotony)</li> <li>• Frequent Itemset Mining (Search space, apriori method)</li> <li>• Association rules (computation, interestingness)</li> <li>• Advanced algorithms for frequent itemset computation [optional]</li> <li>• Data structures to facilitate frequent itemset mining</li> </ul>
Prüfung	<p>Prüfungsleistung (benotet): -Klausur: 90 min</p> <p>Wiederholungsklausur zu Ende des Semesters. Details werden zu Beginn des Moduls bekannt gegeben.</p> <p>In der Klausur weisen die Studierenden nach, inwieweit sie die vorgestellten Prozesse, Modelle und Verfahren verstanden haben, komprimiert wiedergeben und anwenden sowie auf verwandte Problemstellungen übertragen können. In der Klausur werden 7 bis 10 Aufgaben gestellt, die eine eigenständige Beschreibung der Prozessdefinition, Auswahl und Anwendung grundlegender Verfahren (wie zum Beispiel die Anwendung von Klassifikatoren, Clustering Verfahren und Methoden des Frequent Itemset Mining), Evaluierung von Ergebnissen und Entwurf eines Knowledge Discovery Prozesses erfordern.</p>
Examination	<p>Examination requirements (graded): - written exam: 90 min</p> <p>A makeup exam will be offered at the end of the semester, details will be announced at the beginning of the course.</p> <p>Within the written exam, students demonstrate that they understand the presented processes, models, and</p>

	<p>methods, that they can reproduce and apply them as well as that they can transfer and extend models and methods to similar problems. The written exam consist of 7 to 9 assignments, which require the description of process definitions, selection and application of the basic methods (e.g. the application of classifiers, clustering methods and method for frequent itemset mining etc.), the evaluation of the results and the development of a KDD process for a given task.</p>
Literatur/Literature	<p>Han J., Kamber M., Pei J. Data Mining: Concepts and Techniques, 3. Auflage, Morgan Kaufmann, 2011  Tan P.-N., Steinbach M., Kumar V.: Introduction to Data Mining, Addison-Wesley, 2006  Mitchell T. M.: Machine Learning, McGraw-Hill, 1997  Ester M., Sander J.: Knowledge Discovery in Databases: Techniken und Anwendungen, Springer Verlag, 2000  Witten I. H., Frank E., Hall M. A.: Data Mining: Practical Machine Learning Tools and Techniques 3. Auflage, Morgan Kaufmann, 2011</p>
Medienformen	Beamer-Präsentation, Tafelpräsentation, Handout
Media	slides show, blackboard presentation, handouts
Lehr- und Lernmethode	<p>Vorlesung, Übung, Aufgaben zum Selbststudium. Das Modul besteht aus einer Vorlesung und Übungen in kleinen Gruppen (als Tutorübungen). In den Hausaufgaben, die freiwillig abzugeben sind, analysieren die Studierenden die in der Vorlesung vorgestellten Prozesse, Modelle und Verfahren, wenden diese auf konkrete Daten an und erweitern diese für ähnliche Problemstellungen. In den Hausaufgaben werden selbständig anspruchsvolle Übungsaufgaben bearbeitet, die ähnlich zu den Klausuraufgaben sind (siehe oben) und daher zur Vorbereitung darauf dienen. In der Übung werden mögliche Lösungsstrategien der Aufgaben zum Selbststudium diskutiert.</p>
Teaching and Learning Methods	<p>Lecture, tutorial, assignments for individual study. Within the assignments (the submission is optional) students analyze the processes, models, and methods presented in the corresponding lectures, apply them to real data, and extend these to similar problems. The assignments consist of demanding problems similar to the assignments in the written exam (for details see above) and serve as a preparation for the exam. Within the tutorials possible approaches for solutions of the assignments will be discussed.</p>
Turnus	Sommersemester
Modulverantwortlicher	PD Dr. Matthias Schubert

Dozenten	PD Dr. Matthias Schubert
----------	--------------------------