



## Principali informazioni sull'insegnamento

Denominazione dell'insegnamento	<b>Data Mining</b>	
Corso di studio	Laurea Magistrale in Data Science	
Anno Accademico	2024/25	
Crediti formativi universitari (CFU) / European Credit Transfer and Accumulation System (ECTS)	9 CFU	
Settore Scientifico Disciplinare	IINF-05/A - Sistemi di elaborazione delle informazioni	
Lingua di erogazione	Italiano	
Anno di corso	Primo	
Periodo di erogazione	2^ semestre, le date esatte sono riportate nel manifesto/regolamento	
Obbligo di frequenza	No, ma la frequenza è fortemente raccomandata	
Sito web del corso di studio	<a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea</a>	

Docente/i	
Nome e cognome	Donato Malerba
Indirizzo mail	donato.malerba@uniba.it
Telefono	080 5443269
Sede	Dipartimento di Informatica, Via Orabona 4, 70125, Bari. Stanza n.508, 5^ piano.
Sede virtuale	Piattaforma e-learning UNIBA - <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a>
Sito web del docente	<a href="https://www.uniba.it/it/docenti/malerba-donato">https://www.uniba.it/it/docenti/malerba-donato</a>
Ricevimento (giorni, orari e modalità, es. su appuntamento)	Mercoledì 11.00-13.00 o su appuntamento

Syllabus	
Obiettivi formativi	Acquisizione di adeguate conoscenze dei processi di scoperta della conoscenza nei dati (KDD) e di fondamentali tecniche di data mining per analisi di raggruppamento



	<p>e di associazione. Acquisizione di competenza nell'analisi di un dataset con strumenti di data mining, specificatamente nell'utilizzo di algoritmi di clustering e scoperta di regole di associazione.</p>
<b>Prerequisiti</b>	<p>Le seguenti conoscenze preliminari facilitano ed accelerano la comprensione degli argomenti dell'insegnamento:</p> <p>Da <b>FONDAMENTI DI MATEMATICA</b>: nozioni di insiemistica, numeri interi e reali, funzioni reali di variabile reale (valore assoluto, logaritmo, esponenziale, funzioni trigonometriche), limiti e continuità per funzioni di una variabile, derivate di una funzione in una variabile, integrali di funzioni in una variabile, cenni di combinatorica, calcolo matriciale, autovalori e autovettori, relazioni, relazioni funzionali, di equivalenza e di ordine. Reticolo.</p> <p>Da <b>FONDAMENTI DI PROGRAMMAZIONE</b>: algoritmi, linguaggi di programmazione e programmi, flusso di controllo, funzioni, strutture dati, stringhe e file, ricorsione, ricerca, ordinamento, complessità computazionale. Python.</p> <p>Da <b>GESTIONE DI DATI STRUTTURATI E NON STRUTTURATI</b>: Dati, Informazioni, Conoscenza. Dati strutturati e non strutturati. Base di dati e Sistema Informatico. Transazioni. Modello Concettuale e sua rappresentazione in diagrammi E/R. Modello logico relazionale. Fondamenti di SQL. Rappresentazione e interrogazione di dati spaziali, temporali e testuali.</p> <p>Da <b>MODELLIZZAZIONE STATISTICA</b>: distribuzione di probabilità, principali distribuzioni univariate. Probabilità condizionata.</p>
<b>Contenuti di insegnamento (Programma)</b>	<ol style="list-style-type: none"><li>1. Scoperta di conoscenza nelle basi di dati: il processo. (36 ore con laboratorio) La scoperta di conoscenza nelle basi di dati: definizione. Il processo della scoperta di conoscenza nelle basi di dati. Il processo CRISP-DM: business understanding, data understanding, data preparation, modelling, evaluation, deployment.</li><li>2. Similarità e distanze. (18 ore con laboratorio) Similarità/distanze per dati qualitativi e quantitativi, misure di similarità di testi, misure di similarità temporale.</li><li>3. Analisi di associazione. (18 ore con laboratorio) insieme frequente, regola di associazione, principali esempi di pattern di associazione. Applicazioni all'analisi di cestino, di log e alla bioinformatica.</li><li>4. Analisi di raggruppamento (clustering). (14 ore con laboratorio) Algoritmi di clustering partizionale, gerarchico, basati su modelli probabilistici, basati su griglia e densità. Validazione di cluster. Applicazioni a sistemi di raccomandazione, analisi di reti sociali, marketing, biologia e pianificazione urbana.</li></ol>
<b>Testi di riferimento</b>	<p>Charu C. Aggarwal <i>Data Mining</i> Springer 2015 (disponibile in biblioteca e su piattaforma Ada)</p> <p>Gli studenti che lo desiderano possono ottenere i testi in prestito dalla Biblioteca. Può convenire verificarne la disponibilità mediante il Sistema Bibliotecario di Ateneo <a href="https://opac.uniba.it/easyweb/w8018/index.php?">https://opac.uniba.it/easyweb/w8018/index.php?</a> e contattare la biblioteca per concordare il prestito.</p>
<b>Note ai testi di riferimento</b>	<p>I testi di riferimento sono supportati da articoli scientifici e dispense forniti dal docente durante lo svolgimento del corso e disponibili sulla piattaforma e-learning di UNIBA.</p>



					Sulla stessa piattaforma sono resi disponibili esempi di analisi di dataset, richiesta per la discussione orale, e tracce di prove scritte/esonero.
<b>Organizzazione della didattica</b>					
<b>Ore</b>					
Totali	Didattica frontale	Laboratorio/esercitazioni	Progetto	Studio individuale	
225 ore	56 ore	30 ore	0 ore	139 ore	
<b>CFU/ETCS</b>					
9 CFU	7 CFU	2 CFU	0 CFU		

<b>Metodi didattici</b>		
		Le 86 ore previste per l'insegnamento sono ripartite come segue: 56 ore di didattica frontale (in presenza); 30 ore di laboratorio (in presenza).

<b>Risultati di apprendimento previsti</b>		
<b>Conoscenza e capacità di comprensione</b>		<ul style="list-style-type: none"> <li>○ Acquisizione di conoscenze relative agli algoritmi di data mining più noti in letteratura.</li> <li>○ Comprensione delle scelte di algoritmi di data mining per specifici compiti.</li> <li>○ Capacità di interpretazione dei risultati di un algoritmo di data mining.</li> </ul>
<b>Conoscenza e capacità di comprensione applicate</b>		<ul style="list-style-type: none"> <li>○ Capacità di realizzazione di un semplice progetto di scoperta di conoscenza in una collezione di dati mediante: <ul style="list-style-type: none"> <li>▪ Utilizzo e/o sviluppo di strumenti per la selezione, pre-elaborazione e trasformazione dei dati, e per la validazione dei pattern estratti.</li> <li>▪ Utilizzo di strumenti di data mining per l'estrazione di conoscenza finalizzata a scopi descrittivi in diversi contesti applicativi (aziendali e scientifici).</li> </ul> </li> </ul>
<b>Competenze trasversali</b>		<p><b>Autonomia di giudizio</b></p> <ul style="list-style-type: none"> <li>○ Gli studenti sono in grado di apprezzare l'uso di algoritmi di data mining in processi di scoperta della conoscenza.</li> <li>○ L'autonomia di giudizio viene acquisita attraverso lo studio e l'interpretazione critica dei testi.</li> <li>○ Il raggiungimento dell'adeguata autonomia è verificato attraverso le esercitazioni, che si tengono durante il corso, e con l'esame finale di profitto.</li> </ul> <p><b>Abilità comunicative</b></p> <ul style="list-style-type: none"> <li>○ Gli studenti sono in grado di esporre le tematiche incluse nel programma del corso mediante il lessico specifico della disciplina, e sono in grado di comunicare</li> </ul>



	<p>in modo chiaro e privo di ambiguità le loro conclusioni nell'analisi di un dataset, a interlocutori specialisti e non specialisti.</p> <p><b>Capacità di apprendere in modo autonomo</b></p> <ul style="list-style-type: none"><li>○ Gli studenti sono in grado di approfondire in autonomia le tematiche incluse nel programma del corso anche ricorrendo a risorse (articoli su algoritmi di data mining e/o strumenti di analisi) non direttamente utilizzate nelle ore di lezione/laboratorio.</li></ul>
--	---

Valutazione	
<b>Modalità di verifica dell'apprendimento</b>	<ul style="list-style-type: none"><li>○ Prova scritta sulla parte teorica.</li><li>○ Discussione dell'analisi di un dataset assegnato dal docente, mediante l'applicazione della metodologia CRISP-DM vista nel corso e di strumenti scelti dallo studente, al fine di scoprire pattern e cluster di interesse.</li></ul> <p>La prova scritta è propedeutica alla discussione dell'analisi del dataset. La votazione finale è in trentesimi. Ad essa contribuiscono il voto della prova scritta e la valutazione della discussione dell'analisi del dataset.</p> <p>Sono previste due prove parziali, una a metà corso e una a fine corso, che nell'insieme sostituiscono la prova scritta sulla parte teorica. Le prove parziali sono mantenute per il solo anno accademico in cui sono state sostenute.</p> <p>Durante la prova scritta è possibile utilizzare una calcolatrice, mentre nella discussione dell'analisi del dataset è possibile utilizzare il proprio PC dove sono installati gli ambienti utilizzati per l'analisi del dataset.</p>
Criteria di valutazione	<p><b>Conoscenza e capacità di comprensione:</b></p> <ul style="list-style-type: none"><li>○ esposizione critica dei concetti appresi relativi al processo di scoperta della conoscenza e capacità di affrontare semplici esercizi di data mining.</li></ul> <p><b>Conoscenza e capacità di comprensione applicate:</b></p> <ul style="list-style-type: none"><li>○ analisi di dataset con applicazione di algoritmi di data mining e comprensione dei risultati ottenuti, in un ciclo finalizzato al miglioramento delle prestazioni</li></ul> <p><b>Autonomia di giudizio:</b></p> <ul style="list-style-type: none"><li>○ Capacità di svolgere semplici esercizi assegnati durante il corso delle lezioni</li></ul> <p><b>Abilità comunicative:</b></p> <ul style="list-style-type: none"><li>○ Uso del lessico specifico della disciplina informatica</li></ul> <p><b>Capacità di apprendere:</b></p> <ul style="list-style-type: none"><li>○ sviluppo di argomenti su data mining non direttamente trattati nel corso ma assegnati dal docente</li></ul>
Criteria di misurazione dell'apprendimento e di attribuzione del voto finale	<p>Il voto finale sarà attribuito sulla base della valutazione comparativa di tutti i criteri di valutazione sopra citati</p>
<b>Altro</b>	<p>Si suggerisce agli studenti di affidarsi esclusivamente alle informazioni/comunicazioni fornite sui siti ufficiali del Dipartimento di Informatica, ovvero sui gruppi social solo se costituiti e amministrati esclusivamente dai docenti dei relativi insegnamenti:</p> <ul style="list-style-type: none"><li>● <a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea</a></li><li>● <a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica">https://www.uniba.it/it/ricerca/dipartimenti/informatica</a></li></ul>



- <https://elearning.uniba.it/>

I programmi di tutti gli insegnamenti sono disponibili al seguente link:

- <https://elearning.uniba.it/course/index.php?categoryid=456>

Le informazioni che tutti gli studenti dovrebbero conoscere sono scritte nei regolamenti didattici dei Corsi di Studi disponibili nel sito:

- <https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea>

Si suggerisce agli studenti di diffidare delle informazioni e dei materiali circolanti su siti o gruppi social non ufficiali, poiché spesso sono risultati non affidabili, non corretti o incompleti. Per ogni dubbio, chiedere un incontro al docente secondo le modalità previste per il ricevimento.

Link al corso sulla piattaforma e-learning ADA:

<https://elearning.uniba.it/course/index.php?categoryid=456>



## Main information on the course

Course name	<b>Data Mining</b>	
Degree	Master's Degree in Data Science	
Academic year	2024/25	
European Credit Transfer and Accumulation System (ECTS), in Italian Crediti Formativi Universitari (CFU)	9 CFU (each CFU corresponds to 25 hours (h) of student's time); CFU are of type T1, T2 or T3 T1 = 8 h lecture + 17 h individual study T2 = 15 h practice + 10 h individual study T3 = 25 h individual study	
Settore Scientifico Disciplinare		
Course language	Italian	
Course year	First	
Course period	Second Semester - exact dates can be found in the didactic regulations	
Course attendance requirement	None, but it is highly recommended to attend classes	
Website of the Degree	<a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea</a>	

## Teacher(s)

Name and Surname	Donato Malerba
email	donato.malerba@uniba.it
phone	080 5443269
office	Dipartimento di Informatica, Via Orabona 4, 70125, Bari. Room n.508, 5th floor
e-learning platform	Platform ADA - <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a>
Teacher's homepage	<a href="https://www.uniba.it/it/docenti/malerba-donato">https://www.uniba.it/it/docenti/malerba-donato</a>
Office hours	Wednesday 11:00-13:00 or on appointment

## Syllabus

<b>Course goals</b>	Acquisition of adequate knowledge of Knowledge Discovery in Databases (KDD) processes and fundamental data mining techniques for cluster analysis and association analysis. Gaining proficiency in analyzing a dataset using data mining tools, specifically in the utilization of clustering algorithms and association rule discovery.
<b>Prerequisites/requirements</b>	<p>The following preliminary knowledge facilitates and accelerates the understanding of the course topics:</p> <p>From FUNDAMENTALS OF MATHEMATICS: Notions of set theory, integers and real numbers, real functions of a real variable (absolute value, logarithm, exponential, trigonometric functions), limits and continuity for functions of one variable, derivatives of a function in one variable, integrals of functions in one variable, basics of combinatorics, matrix calculus, eigenvalues and eigenvectors, relations, functional relations, equivalence relations, and order relations. Lattice.</p> <p>From FUNDAMENTALS OF PROGRAMMING: Algorithms, programming languages, and programs, flow control, functions, data structures, strings, and files, recursion, search, sorting, computational complexity. Python.</p> <p>From MANAGEMENT OF STRUCTURED AND UNSTRUCTURED DATA:</p>



	<p>Data, Information, Knowledge. Structured and unstructured data. Databases and Information Systems. Transactions. Conceptual model and its representation in E/R diagrams. Logical relational model. Basics of SQL. Representation and querying of spatial, temporal, and textual data.</p> <p>From STATISTICAL MODELING: Probability distribution, main univariate distributions. Conditional probability.</p>			
<b>Course program</b>	<p>1. Knowledge Discovery in Databases (KDD): The Process. (36 hours including practice sessions) Knowledge Discovery in Databases: Definition. The process of Knowledge Discovery in Databases. The CRISP-DM process: business understanding, data understanding, data preparation, modeling, evaluation, deployment.</p> <p>2. Similarity and Distances. (18 hours including practice sessions) Similarity/distances for qualitative and quantitative data, measures of text similarity, measures of temporal similarity.</p> <p>3. Association Analysis. (18 hours including practice sessions) Frequent itemsets, association rules, main examples of association patterns. Applications in market basket analysis, log analysis, and bioinformatics.</p> <p>4. Clustering Analysis. (14 hours including practice sessions) Partitioning, hierarchical, probabilistic model-based, grid-based, and density-based clustering algorithms. Cluster validation. Applications in recommendation systems, social network analysis, marketing, biology, and urban planning.</p>			
<b>Books of reference</b>	<p>Charu C. Aggarwal <i>Data Mining</i> Springer 2015 (available in the library and on the Ada platform)</p> <p>Students who wish to can borrow the texts from the Library. It may be advisable to check their availability through the University Library System at <a href="https://opac.uniba.it/easyweb/w8018/index.php?">https://opac.uniba.it/easyweb/w8018/index.php?</a> and contact the library to arrange the loan.</p>			
<b>Notes to the books</b>	<p>The reference texts are supported by scientific articles and slides provided by the instructor during the course.</p>			
<b>Organization of the didactic activities</b>				
<b>Hours</b>				
Total	Lectures	Practice sessions	Project work	Individual study
225 hours	56 hours	30 hours	0 hours	139 hours
<b>CFU/ETCS</b>				
9 CFU	7 CFU	2 CFU	0 CFU	
<b>Teaching methods</b>				
	<p>The 86 hours allocated for classes are distributed as follows:</p> <ul style="list-style-type: none"> <li>- 56 hours of lectures (in person);</li> <li>- 30 hours of practice sessions (in person).</li> </ul>			
<b>Expected learning outcomes</b>				



<b>Knowledge and understanding</b>	<ul style="list-style-type: none"> <li>○ Acquisition of knowledge related to the most well-known data mining algorithms in the literature.</li> <li>○ Understanding the choices of data mining algorithms for specific tasks.</li> <li>○ Ability to interpret the results of a data mining algorithm.</li> </ul>
<b>Applying knowledge and understanding</b>	<ul style="list-style-type: none"> <li>○ Ability to carry out a simple knowledge discovery project in a data collection through:             <ul style="list-style-type: none"> <li>▪ Utilization and/or development of tools for data selection, pre-processing, transformation, and validation of extracted patterns.</li> <li>▪ Usage of data mining tools for knowledge extraction aimed at descriptive purposes in various application contexts (business and scientific).</li> </ul> </li> </ul>
<b>Other skills</b>	<p><i>Making judgements</i></p> <ul style="list-style-type: none"> <li>○ Students are able to appreciate the use of data mining algorithms in knowledge discovery processes.</li> <li>○ Independent judgment is acquired through the study and critical interpretation of texts.</li> <li>○ The achievement of adequate autonomy is assessed through exercises held during the course and the final exam.</li> </ul> <p><i>Communication</i></p> <ul style="list-style-type: none"> <li>○ Students are able to present the topics included in the course syllabus using the specific terminology of the discipline, and they are capable of clearly and unambiguously communicating their conclusions in the analysis of a dataset to both specialist and non-specialist audiences.</li> </ul> <p><i>Learning skills</i></p> <ul style="list-style-type: none"> <li>○ Students are able to independently delve deeper into the topics covered in the course syllabus, even by using resources (such as articles on data mining algorithms and/or analysis tools) that are not directly explained during lectures or practice sessions.</li> </ul>

<b>Assessment</b>	
<b>Assessment methods</b>	<ul style="list-style-type: none"> <li>○ Written Exam on Theoretical Content: A written exam covering the theoretical part of the course.</li> <li>○ Discussion of a Dataset Analysis: Students must discuss the analysis of a dataset assigned by the instructor, applying the CRISP-DM methodology covered during the course and using tools chosen by the student. The goal is to discover patterns and clusters of interest.</li> </ul> <p>The written exam is a prerequisite for the discussion of the dataset analysis. The final grade is a vote between 0 and 30 (cum laude) and is based on the score of the written exam and the evaluation of the dataset analysis discussion.</p> <p>Two partial exams are offered, one at mid-course and one at the end of the course, which together replace the written exam on the theoretical content. These partial exams are valid only for the academic year in which they are taken.</p> <p>During the written exam, students are allowed to use a pocket calculator, while during the dataset analysis discussion, they can use their own PC with the necessary software environments installed for dataset analysis.</p>





<b>Evaluation criteria</b>	<p><b>Knowledge and Understanding:</b></p> <ul style="list-style-type: none"><li>○ Critical presentation of the concepts learned related to the knowledge discovery process and the ability to tackle simple data mining exercises.</li></ul> <p><b>Applied Knowledge and Understanding:</b></p> <ul style="list-style-type: none"><li>○ Analysis of datasets with the application of data mining algorithms and understanding of the obtained results, in a cycle aimed at improving performance.</li></ul> <p><b>Autonomy of Judgment:</b></p> <ul style="list-style-type: none"><li>○ Ability to carry out simple exercises assigned during the course lectures.</li></ul> <p><b>Communication Skills:</b></p> <ul style="list-style-type: none"><li>○ Use of specific vocabulary in the field of computer science.</li></ul> <p><b>Ability to Learn:</b></p> <ul style="list-style-type: none"><li>○ Development of topics on data mining not directly covered in the course but assigned by the instructor.</li></ul>
Measurements and final grade	The final grade will be awarded based on the comparative evaluation of all the assessment criteria mentioned above.
<b>Further information</b>	<p>It is recommended that students rely exclusively on information and communications provided on the official websites of the Department of Computer Science or on social groups only if they are established and managed solely by the instructors of the respective courses:</p> <ul style="list-style-type: none"><li>● <a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea</a></li><li>● <a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica">https://www.uniba.it/it/ricerca/dipartimenti/informatica</a></li><li>● <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a></li></ul> <p>The course syllabi are available here:</p> <ul style="list-style-type: none"><li>● <a href="https://elearning.uniba.it/course/index.php?categoryid=456">https://elearning.uniba.it/course/index.php?categoryid=456</a></li></ul> <p>The information that all students should know is detailed in the educational regulations and course handbooks available on the website:</p> <ul style="list-style-type: none"><li>● <a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/corsi-di-laurea</a></li></ul> <p>Students are advised to be cautious regarding information and materials circulating on unofficial websites or social groups, as they can often be unreliable, incorrect, or incomplete. For any doubts, students should arrange a meeting with the instructor following the specified office hours.</p> <p>Link to the course on the ADA e-learning platform: <a href="https://elearning.uniba.it/course/index.php?categoryid=456">https://elearning.uniba.it/course/index.php?categoryid=456</a></p>