



## Principali informazioni sull'insegnamento

Denominazione dell'insegnamento	<b>Big Data</b>	
Corso di studio	Computer Science (Masters' degree in Computer Science)	
Anno Accademico	2023/24	
Crediti formativi universitari (CFU) / European Credit Transfer and Accumulation System (ECTS)	6 CFU	
Settore Scientifico Disciplinare	INF/01	
Lingua di erogazione	inglese	
Anno di corso	Secondo	
Periodo di erogazione	1^ semestre, le date esatte sono riportate nel manifesto/regolamento	
Obbligo di frequenza	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/computer-science/computer-science">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/computer-science/computer-science</a>	

<b>Docente/i</b>	
Nome e cognome	Michelangelo Ceci
Indirizzo mail	michelangelo.ceci@uniba.it
Telefono	+39 080 5442285
Sede	Dipartimento di Informatica, Via Orabona 4, 70125, Bari. Stanza n.510, 5^ piano.
Sede virtuale	Piattaforma ADA - <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a>
Sito web del docente	<a href="http://www.di.uniba.it/~ceci/">http://www.di.uniba.it/~ceci/</a>
Ricevimento (giorni, orari e modalità, es. su appuntamento)	In presenza presso lo studio del docente: Giovedì, 9:00-11:00, preferibilmente su appuntamento

## Syllabus



<b>Obiettivi formativi</b>	<p>Conoscenza e comprensione</p> <p>L'analisi dei dati ha sostituito l'acquisizione dei dati come collo di bottiglia per il processo decisionale basato sull'evidenza. Ciò implica che l'estrazione di conoscenza da set di dati grandi, eterogenei e rumorosi richiede non solo potenti risorse di calcolo, ma anche basi metodologiche e astrazioni di programmazione appropriate. Le astrazioni emerse nell'ultimo decennio fondono idee da database paralleli, sistemi distribuiti e linguaggi di programmazione per creare una nuova classe di piattaforme di analisi dei dati scalabili che costituiscono la base per la scienza dei dati su scala reale. Nel corso verrà fornita una panoramica dei principi e tecnologie rilevanti, sia in termini di organizzazione e gestione dei dati che in termini di analisi degli stessi.</p> <p>Conoscenza e comprensione applicate. Progettazione di un processo di analisi di grandi volumi di dati, considerando la loro gestione e memorizzazione, la loro elaborazione e la loro analisi.</p>
<b>Prerequisiti</b>	È necessario conoscere i principi base delle basi di dati, le basi di dati distribuite, i Data Warehouse, i concetti basilari di Apprendimento Automatico.
<b>Contenuti di insegnamento (Programma)</b>	<p><b>1. Big Data: introduction (5h).</b> Big Data: difinitions, Why and where, Characteristics of Big Data dimensions and scalability, Getting value from Big Data, A fourth paradigm for doing science.</p> <p><b>2. Big data modeling and storing (5h).</b> The aggregate model. The key-value, document and column family models. L'esempio di BigTable. Column-family: Google BigTable; Dynamo DB; Graph databases: concepts, queries and data organization</p> <p><b>3. Knowledge Discovery in Databases (10h).</b> Context and definitions. The KDD process. Crisp_DM: Business understanding, Data Understanding, Sampling. Feature selection, Feature transformation and construction. Factor analysis. Data representation, Data representation, model representation, Model building and evaluation; PMML.</p> <p><b>4. Correlation and variable associations (4h).</b> Correlation Between two/many quantitative variables. Correlation between qualitative variables: Chi-square, lambda, Sparman Correlation Coefficient. Association rule mining: Apriori and FPGrowth.</p> <p><b>5. The Bayesian framework and Ensemble learning (5h).</b> MAP learning, Bayes optimal classifier. Ensemble learning: Boosting, Bagging and Stacking</p> <p><b>6. Programming framework for Big Data Analysis (5h).</b> Spark: The MapReduce framework, RDDs, DataFrames and DataBases. Spark SQL, Spark MLlib, Park streaming.</p> <p><b>7. Learning from data streams (4h).</b> Motivations, Data Streams definitions, Approximate Answers, Count-Min Sketch. Basic Methods: Estimating statistics over windows, Sampling. Very Fast Decision Trees, Conceptd Drift</p> <p><b>Laboratorio (30h)</b> Cassandra DB: Introduction, the architecture of Cassandra, The logical Model. Hand-on sessions Spark, SparkML and Spark streaming., Hand-on sessions.</p>
<b>Testi di riferimento</b>	Viktor Mayer-Schonberger, Kenneth Cukier. Big Data: A Revolution That Will Transform How We Live, Work, and Think, John Murray, 2013



	<p>T. Mitchell Machine Learning, Morgan Kaufmann, 1997 Richard J. Roiger, Michael W. Geatz. Introduction to Data Mining McGraw-Hill, 2003 A. Azzalini, B. Scarpa Analisi dei dati e data mining, Springer, 2004</p> <p>Articoli scientifici selezionati e messi a disposizione su ADA. Slide presentate a lezione e messe a disposizione su 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>Sul sito ADA sono disponibili: Le slide presentate a lezione, slide presentate durante i seminari e articoli scientifici di interesse.</p>		
<b>Organizzazione della didattica</b>			
<b>Ore</b>			
Totali	Didattica frontale	Pratica (laboratorio, progetto, esercitazione, altro)	Studio individuale
150 ore	32 ore	30 ore	88 ore
<b>CFU/ETCS</b>			
6 CFU	4 CFU	2 CFU	

<b>Metodi didattici</b>	
	<p>Lezioni frontali sugli argomenti teorici del programma e su esempi pratici da sviluppare in aula e a casa. Attività in laboratorio per approfondire gli aspetti pratici e tecnologici.</p>

<b>Risultati di apprendimento previsti</b>	
<b>Conoscenza e capacità di comprensione</b>	<p>Lo studente acquisirà le principali conoscenze riguardanti il panorama dei sistemi rilevanti di Big Data, i principi su cui si basano i metodi di memorizzazione e il processo di analisi basato su metodologia CRISP-DM. Verranno trattati i database NoSQL, Spark e l'ecosistema che esso ha generato.</p>
<b>Conoscenza e capacità di comprensione applicate</b>	<p>Lo studente acquisirà le conoscenze pratiche che gli consentiranno di applicare le conoscenze sui sistemi di data storage, processing e analisi su grandi moli di dati.</p>



<b>Competenze trasversali</b>	<p><b>Autonomia di giudizio</b></p> <p>Effettuare valutazioni e scelte informate è esattamente lo scopo del corso, gli studenti impareranno a progettare e implementare in modo autonomo uno strumento di Big Data Analytics ed eseguire analisi.</p> <p><b>Abilità comunicative</b></p> <p>Per rendere fruibile e credibile la conoscenza estratta, il risultato dell'analisi deve essere presentato adeguatamente. Inoltre, il processo deve essere chiaro basato su standard e trasparente. Questo è un passaggio fondamentale del processo KDD e, di conseguenza, questo è un passaggio fondamentale della Big Data Analytics.</p> <p><b>Capacità di apprendere in modo autonomo</b></p> <p>Lo studente apprenderà i concetti di base che lo renderanno in grado di utilizzare, comprendere e implementare qualsiasi metodo di data mining che estrae conoscenza da un grande volume di dati.</p>
-------------------------------	-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<b>Valutazione</b>	
<b>Modalità di verifica dell'apprendimento</b>	<p>L'esame si compone di una parte orale e di una parte pratica. La parte orale ha lo scopo di verificare i concetti appresi e si riferisce a tutti gli argomenti discussi durante le lezioni. La parte pratica è una discussione orale di un piccolo elaborato pratico e mira a verificare le capacità di creare un processo CRISP-DM per l'analisi dei big data.</p> <p>Le due parti non sono propedeutiche l'una all'altra (cioè sono indipendenti l'una dall'altra). I voti ottenuti per entrambe le parti saranno mantenuti fino al termine dell'Anno Accademico.</p> <p>La valutazione è in trentesimi per entrambe le prove e il voto finale è derivato applicando la media delle valutazioni delle prove.</p>
Criteri di valutazione	<p>La valutazione mira a verificare il raggiungimento di una buona capacità di problem solving nell'ambito dell'analisi di grandi volumi di dati. In particolare, rispetto ai risultati di apprendimento attesi, si considerano i seguenti criteri:</p> <p>Conoscenza e capacità di comprensione: Capacità di descrivere con accuratezza e chiarezza le tecniche e i fondamenti teorici del processo di analisi dei Big Data;</p> <p>Conoscenza e capacità di comprensione applicate: Capacità di progettare e definire una architettura per l'analisi dei dati.</p> <p>Autonomia di giudizio: Capacità di giustificare la soluzione adottata tra più opzioni;</p> <p>Abilità comunicative: Chiarezza di risposte/soluzioni a domande/problemi;</p> <p>Capacità di apprendere:</p>



	<p>Astrazione, ragionamento per analogia e dimostrazione di adattività nella risoluzione dei problemi, anche considerando il dominio di riferimento e i dati a disposizione del processo di analisi.</p>
<p>Criteria di misurazione dell'apprendimento e di attribuzione del voto finale</p>	<p>Lo studente dovrà dimostrare di conoscere tutti i concetti discussi durante le lezioni, nonché dimostrare di essere in grado di progettare una analisi completa secondo il modello CRISP-DM, realizzato mediante parte delle tecnologie discusse.</p>
<p><b>Altro</b></p>	<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/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>I programmi degli insegnamenti sono disponibili qui:</p> <ul style="list-style-type: none"><li>• <a href="https://elearning.uniba.it/course/index.php?categoryid=287">https://elearning.uniba.it/course/index.php?categoryid=287</a></li></ul> <p>Le informazioni che tutti gli studenti dovrebbero conoscere sono scritte nei Regolamenti didattici e manifesti degli studi disponibili nel sito:</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>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.</p>



## Main information on the course

Course name	<b>Big Data</b>	
Degree	Computer Science (Masters' degree in Computer Science)	
Academic year	2023/24	
European Credit Transfer and Accumulation System (ECTS), in Italian Crediti Formativi Universitari (CFU)	6 CFU	
Scientific Disciplinary Sector	ING-INF/05	
Course language	English	
Academic year	Second	
Delivery period	First semester, exact dates are specified in the program/regulations.	
Attendance requirement	It is highly recommended to attend classes	
Course of study's website	<a href="https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/computer-science/computer-science">https://www.uniba.it/it/ricerca/dipartimenti/informatica/didattica/corsi-di-laurea/computer-science/computer-science</a>	

Teacher(s)	
Name and Surname	Michelangelo Ceci
Email	michelangelo.ceci@uniba.it
Phone	+39 080 5442285
Office	Dipartimento di Informatica, Via Orabona 4, 70125, Bari. Room n.510, 5th floor.
E-learning platform	Platform ADA - <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a>
Teacher's homepage	<a href="http://www.di.uniba.it/~ceci/">http://www.di.uniba.it/~ceci/</a>
Office hours	In-person at the office: Thursday, 9:00-11:00 a.m. preferably by appointment



Syllabus	
<b>Course goals</b>	<p>Knowledge and Understanding</p> <p>The analysis of data has replaced data acquisition as the bottleneck for evidence-based decision-making processes. This implies that extracting knowledge from large, heterogeneous, and noisy datasets requires not only powerful computing resources but also appropriate methodological foundations and programming abstractions. The abstractions that have emerged in the last decade merge ideas from parallel databases, distributed systems, and programming languages to create a new class of scalable data analytics platforms that form the basis for real-world scale data science. The course will provide an overview of relevant principles and technologies, both in terms of data organization and management, and in terms of data analysis.</p> <p>Applied Knowledge and Understanding</p> <p>Designing a process for analyzing large volumes of data, considering their management and storage, processing, and analysis.</p>
<b>Prerequisites/requirements</b>	<p>It is necessary to be familiar with the basic principles of databases, distributed databases, Data Warehouses, and fundamental concepts of Machine Learning.</p>
<b>Course program</b>	<p><b>Big Data: Introduction (5h):</b> Definitions of Big Data, reasons and contexts for its use, characteristics of Big Data, scalability dimensions, extracting value from Big Data, and the fourth paradigm for scientific research.</p> <p><b>Big Data Modeling and Storing (5h):</b> The aggregate model, key-value, document, and column family models. Examples of BigTable, column-family databases like Google BigTable, Dynamo DB, and concepts, queries, and data organization in graph databases.</p> <p><b>Knowledge Discovery in Databases (10h):</b> Context and definitions, the KDD process, Crisp_DM (Business understanding, Data Understanding, Sampling, Feature selection, Feature transformation, and construction). Factor analysis, data representation, model representation, model building and evaluation, PMML.</p> <p><b>Correlation and Variable Associations (4h):</b> Correlation between two/many quantitative variables, correlation between qualitative variables (Chi-square, lambda, Spearman Correlation Coefficient). Association rule mining using Apriori and FPGrowth.</p> <p><b>The Bayesian Framework and Ensemble Learning (5h):</b> MAP learning, Bayes optimal classifier, ensemble learning techniques including Boosting, Bagging, and Stacking.</p> <p><b>Programming Framework for Big Data Analysis (5h):</b> Introduction to Spark, the MapReduce framework, RDDs, DataFrames, and Databases. Overview of Spark SQL, Spark MLlib, and Spark streaming.</p> <p><b>Learning from Data Streams (4h):</b> Motivations, definitions of data streams, approximate answers, Count-Min Sketch. Basic methods for estimating statistics over windows, sampling, very fast decision trees, and concept drift.</p> <p><b>Laboratory (30h):</b></p> <p><b>Cassandra DB:</b> Introduction, architecture of Cassandra, logical model. Hands-on sessions.</p> <p><b>Spark, SparkML, and Spark streaming:</b> Hands-on sessions.</p>



<b>Books of reference</b>	<p>Viktor Mayer-Schonberger, Kenneth Cukier. Big Data: A Revolution That Will Transform How We Live, Work, and Think, John Murray, 2013  T. Mitchell Machine Learning, Morgan Kaufmann, 1997  Richard J. Roiger, Michael W. Geatz. Introduction to Data Mining McGraw-Hill, 2003  A. Azzalini, B. Scarpa Analisi dei dati e data mining, Springer, 2004</p> <p>Selected scientific articles made available on ADA.  Slides presented during the lecture and provided on ADA.</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>On the ADA website, you can find:  The slides presented in class, slides presented during seminars, and relevant scientific articles.</p>		
<b>Organization of the didactic activities</b>			
<b>Hours</b>			
Total	Lectures	Practice sessions	Individual study
150 hours	32 hours	30 hours	88 hours
<b>CFU/ETCS</b>			
6 CFU	4 CFU	2 CFU	

<b>Teaching methods</b>	
	<p>Frontal lectures on the theoretical topics of the program and practical examples to be developed in the classroom and at home.  Laboratory activities to delve into practical and technological aspects.</p>

<b>Expected learning outcomes</b>	
<b>Knowledge and understanding</b>	<p>The student will acquire the main knowledge regarding the landscape of relevant Big Data systems, the principles underlying storage methods, and the analysis process based on the CRISP-DM methodology. Topics covered include NoSQL databases, Spark, and the ecosystem it has generated.</p>
<b>Applying knowledge and understanding</b>	<p>The student will acquire practical knowledge enabling them to apply their understanding of data storage, processing systems, and large-scale data analysis.</p>





<b>Other skills</b>	<p><b>Judgment autonomy</b></p> <p>Making informed evaluations and choices is precisely the aim of the course; students will learn to independently design and implement a Big Data Analytics tool and conduct analyses.</p> <p><b>Communication skills</b></p> <p>To make the extracted knowledge usable and credible, the results of the analysis must be presented appropriately. Moreover, the process must be clear, based on standards, and transparent. This is a fundamental step in the KDD process and, consequently, a fundamental step in Big Data Analytics.</p> <p><b>Self-learning ability</b></p> <p>The student will learn the basic concepts that will enable them to use, understand, and implement any data mining method that extracts knowledge from a large volume of data.</p>
---------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

<b>Assessment</b>	
<b>Assessment methods</b>	<p>The exam consists of an oral part and a practical part. The oral part aims to assess the understanding of acquired concepts and covers all topics discussed during the lectures. The practical part involves an oral discussion of a small practical project and aims to evaluate the ability to create a CRISP-DM process for big data analysis.</p> <p>The two parts are not interdependent (i.e., they are independent of each other). Grades obtained for both parts will be retained until the end of the academic year.</p> <p>The evaluation is in thirtieths for both components, and the final grade is derived by averaging the scores from both parts.</p>
<b>Evaluation criteria</b>	<p>The evaluation aims to assess the achievement of a good problem-solving ability in the context of large-scale data analysis. In particular, considering the expected learning outcomes, the following criteria are taken into account:</p> <p>Knowledge and understanding:</p> <p>The ability to accurately and clearly describe the techniques and theoretical foundations of the Big Data analysis process.</p> <p>Applied knowledge and understanding:</p> <p>The ability to design and define an architecture for data analysis.</p> <p>Judgment autonomy:</p> <p>The ability to justify the chosen solution among multiple options.</p> <p>Communication skills:</p> <p>Clarity in providing answers/solutions to questions/problems.</p> <p>Learning ability:</p> <p>Abstraction, reasoning by analogy, and demonstration of adaptability in problem-solving, considering the reference domain and the data available for the analysis process.</p>



<b>Measurements and final grade</b>	The student must demonstrate knowledge of all the concepts discussed during the lectures and show the ability to design a comprehensive analysis according to the CRISP-DM model, implemented using some of the discussed technologies.
<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=288">https://elearning.uniba.it/course/index.php?categoryid=288</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> <hr/> <p>Link to the course on the department's ADA e-learning platform: <a href="https://elearning.uniba.it/">https://elearning.uniba.it/</a></p>