

Denominazione Insegnamento/Altra Attività	<b>FONDAMENTI DI DATA SCIENCE E MACHINE LEARNING</b>
SSD	INF/01
CFU	9
Ore	72
Tipologia Attività Formativa - Ambito	Obbligatorie – Discipline Informatiche
Obiettivi Formativi: Risultati d'Apprendimento Previsti e Competenze da Acquisire (Descrittori di Dublino)	<p>L'insegnamento mira a fornire le competenze metodologiche e tecnologiche necessarie per estrarre conoscenza da grossi volumi di dati, mediante tecniche di data profiling, data mining e machine learning, utilizzando opportune strategie di visualizzazione dei risultati. In particolare, si intendono integrare le conoscenze di data management acquisite nell'ambito di altri corsi di basi di dati con competenze utili allo svolgimento della professione del data scientist.</p> <p><b>Conoscenza e Capacità di Comprensione</b></p> <p>Le Principali Conoscenze Acquisite Saranno:</p> <ul style="list-style-type: none"> <li>• Big Data</li> <li>• Data wrangling</li> <li>• Estrazione automatica di dipendenze tra i dati</li> <li>• Data quality e data cleansing</li> <li>• Data integration</li> <li>• Data mining</li> <li>• Mapreduce</li> <li>• Funzioni di similarità</li> <li>• Machine learning</li> <li>• Reti neurali</li> </ul> <p><b>Capacità di Applicare Conoscenza e Comprensione</b></p> <p>Gli studenti acquisiranno la capacità di:</p> <ul style="list-style-type: none"> <li>• Acquisire, organizzare, gestire ed elaborare grandi volumi di dati</li> <li>• Estrarre conoscenza dai dati</li> <li>• Selezionare dati utili</li> <li>• Organizzare un progetto basato su tecniche di machine learning</li> <li>• Comunicare la conoscenza estratta dai dati attraverso diverse forme di rappresentazione, incluso tecniche visuali.</li> </ul>
Prerequisiti	Lo studente deve conoscere i fondamenti di data management, sistemi distribuiti, paradigma ad oggetti ed un linguaggio di programmazione.

<p>Contenuti/Programma</p>	<p>Dopo una panoramica sui nuovi scenari applicativi legati alla gestione di grandi collezioni di dati distribuiti ed eterogenei, incluso le potenzialità di tecnologie capaci di estrarre conoscenza dai dati, il corso si concentrerà sui seguenti argomenti:</p> <p><b>Big data (4 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Problematiche di big data (2 ore di teoria)</li> <li>• Tecnologie di supporto ai big data (2 ore di teoria)</li> </ul> <p><b>Data preparation (12 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Data profiling (4 ore di teoria)</li> <li>• Dipendenze funzionali approssimate e loro utilizzo nel data quality (4 ore di teoria)</li> <li>• Integrazione dati da sorgenti multiple (4 ore di teoria)</li> </ul> <p><b>Estrazione di conoscenza da grandi collezioni di dati (12 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Mapreduce (4 ore di teoria)</li> <li>• Valutazione della similarità (5 ore di teoria)</li> <li>• Introduzione al data mining (2 ore di teoria)</li> <li>• Algoritmo apriori (1 ora di teoria)</li> </ul> <p><b>Machine learning (24 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Concetti introduttivi (4 ore di teoria)</li> <li>• Fasi di un progetto di machine learning (5 ore di teoria)</li> <li>• Modelli di Addestramento (2 ore di teoria)</li> <li>• Classificazione/regressione (3 ore di teoria)</li> <li>• Alberi di decisione (2 ore di teoria)</li> <li>• Ensemble learning and random forest (2 ore di teoria)</li> <li>• Clustering (2 ore di teoria)</li> <li>• Riduzione della dimensionalità (2 ore di teoria)</li> <li>• Support vector machine (2 ore di teoria)</li> </ul> <p><b>Reti Neurali (14 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Introduzione alle reti neurali (2 ore di teoria)</li> <li>• Tensor flow (2 ore di teoria)</li> <li>• Percettroni Multilivello e Reti Neurali Profonde (2 ore di teoria)</li> <li>• Reti convoluzionali (2 ore di teoria)</li> <li>• Reti ricorrenti (4 ore di teoria)</li> <li>• Autoencoder (2 ore di teoria)</li> </ul> <p><b>Strumenti per la data science (6 ore di teoria)</b></p> <ul style="list-style-type: none"> <li>• Il linguaggio python (4 ore frontali)</li> <li>• Weka (2 ore frontali)</li> </ul>
<p>Metodi Didattici</p>	<p>L'insegnamento prevede 66 ore di didattica frontale su argomenti teorici e 6 ore su linguaggi e strumenti applicativi, con l'obiettivo di presentare i concetti e sviluppare capacità di progettare e implementare soluzioni per problematiche di data science e machine learning. Gli argomenti del programma vengono presentati con l'ausilio di presentazioni powerpoint, stimolando discussioni critiche con la classe. Per ogni argomento trattato, vengono illustrati possibili task che possono essere oggetto di un progetto di corso da parte di uno o più studenti. Per quanto riguarda gli strumenti applicativi, oltre all'utilizzo di</p>

	<p>presentazioni powerpoint, nelle quali vengono presentati concetti ed eventuali link a forum, manuali e siti di approfondimento, durante l'orario di ricevimento viene offerta agli studenti la possibilità di chiedere supporto in merito a simulazioni da essi effettuate sul proprio computer, di chiedere chiarimenti e risolvere eventuali problemi tecnici insieme al docente.</p>
Modalità di Verifica dell'Apprendimento	<p>Il raggiungimento degli obiettivi dell'insegnamento è certificato mediante il superamento di un esame con valutazione in trentesimi. L'esame prevede una prova scritta (in alternativa, una prova in itinere a metà corso) ed una prova orale. Inoltre, opzionalmente, gli studenti possono sviluppare un progetto per incrementare il voto ottenuto con le suddette prove. La prova scritta (o quella in itinere) mira ad accertare l'acquisizione dei concetti teorici. La prova orale invece consiste in un colloquio con domande e discussione sui contenuti teorici e metodologici trattati a lezione ed è finalizzata ad accertare la capacità di conoscenza e comprensione, nonché la capacità di esposizione dei concetti. La prova orale rappresenta la prova finale, pertanto essa può essere svolta solo dopo il superamento della prova scritta e, qualora si sia optato anche per lo sviluppo del progetto, anche dopo il completamento e la discussione di quest'ultimo.</p> <p>Il progetto è finalizzato ad accertare la capacità di applicare le conoscenze acquisite in scenari reali. Esso può essere svolto individualmente o in gruppi di massimo 3 persone, scegliendo tra un ventaglio di proposte fatte dal docente. Durante lo svolgimento del progetto, gli studenti potranno interagire con il docente al fine di comunicare gli stati di avanzamento dello stesso e le eventuali criticità emerse, concordando obiettivi e modalità di prosecuzione. Al termine del progetto, gli studenti devono consegnare al docente una tesina contenente la documentazione di progetto ricevendo, dopo qualche settimana, una prima valutazione dello stesso. Quest'ultima potrebbe contenere richieste di integrazione e/o di revisione del lavoro svolto. Pertanto, occorre sottomettere i risultati del progetto con diverse settimane di anticipo rispetto alla data in cui si intende sostenere la prova orale, onde consentire al docente di effettuare la correzione ed al gruppo di progetto di apportare le eventuali correzioni richieste. Al termine del progetto, al gruppo di progetto potrebbe essere richiesto di preparare una tesina ed una presentazione powerpoint della durata di circa 30 minuti.</p> <p>Il voto finale scaturisce, generalmente, dalla media dei voti in trentesimi conseguiti alla prova scritta (alternativamente, quella in itinere) ed a quella orale, con la possibilità di incrementare il punteggio così ottenuto fino a 3 punti, tramite lo sviluppo del progetto.</p>
Testi di Riferimento	<ol style="list-style-type: none"> <li>1. Jure Leskovec, Anand Rajaraman, Jeffrey D. Ullman, Mining Of Massive Datasets", 3<sup>a</sup> Edizione, Cambridge University Press, 2020.</li> <li>2. Aurélien Géron, " Hands-On Machine Learning with Scikit-Learn, Keras, and Tensorflow: Concepts, Tools, and Techniques to Build Intelligent Systems “, 3<sup>a</sup> Edizione, O Reilly ed, 2023.</li> <li>3. Chirag Shah, A Hands-On Introduction To Data Science, Cambridge University Press, 2020.</li> </ol>

	<ol style="list-style-type: none"> <li>4. Foster Provost, Tom Fawcett, Data Science For Business: What You Need To Know About Data Mining And Data-Analytic Thinking, O Reilly Ed, 2013.</li> <li>5. P. Deitel, H. Deitel, Introduzione A Python – Per L'informatica E La Data Science, Pearson 2021.</li> </ol>
Altre Informazioni	<p>La frequenza del corso è fortemente consigliata. Gli studenti devono essere preparati a trascorrere una congrua quantità di tempo nello studio al di fuori delle lezioni. Una preparazione soddisfacente richiede in media 1 ora di studio per ciascuna ora trascorsa in aula e circa 80 ore per lo sviluppo del progetto.</p> <p>Il materiale delle lezioni sarà disponibile sulla piattaforma e-learning dipartimentale <a href="http://elearning.informatica.unisa.it/el-platform/">http://elearning.informatica.unisa.it/el-platform/</a>.</p> <p>Contatti</p> <p>Prof. Giuseppe Polese  <a href="mailto:gpolese@unisa.it">gpolese@unisa.it</a></p>