

<b>Principali informazioni sull'insegnamento</b>	
Anno di corso	2023/2024
Periodo di erogazione	marzo 2024 – maggio 2024
Crediti formativi universitari (CFU/ETCS):	6
SSD	ING-INF/05
Lingua di erogazione	Italiano o Inglese
Modalità di frequenza	Frequenza facoltativa

<b>Docente</b>	
Nome e cognome	Dott. Roberto Cilli
Indirizzo mail	<a href="mailto:roberto.cilli@uniba.it">roberto.cilli@uniba.it</a>
Telefono	327 730 8761
Sede	Dipartimento di Fisica, Primo Piano, Stanza 124
Sede virtuale	Canale Teams
Ricevimento	Ogni venerdì – ore 15-19

<b>Organizzazione della didattica</b>			
<b>Ore</b>			
Totali	Didattica frontale	Pratica (laboratorio, campo, esercitazione, altro)	Studio individuale
150	32	30	88
<b>CFU/ETCS</b>			
6	4	2	

<b>Obiettivi formativi</b>	<i>Il corso di insegnamento in machine learning ed intelligenza artificiale ha come obiettivo formativo quello di fornire agli studenti le competenze necessarie per comprendere e implementare algoritmi di apprendimento automatico per l'estrazione di conoscenza in grandi moli di dati multi-omici.</i>
<b>Prerequisiti</b>	<i>Non vi sono prerequisiti specifici differenti da quelli richiesti per l'accesso al corso di laurea. Tuttavia, conoscenze preliminari di programmazione in linguaggio R e di statistica possono aiutare ad affrontare i contenuti previsti dall'insegnamento.</i>

<p><b>Metodi didattici</b></p>	<p><i>La modalità di erogazione dell'insegnamento prevede delle lezioni frontali alternati ad esercitazioni in aula. L'erogazione dell'insegnamento segue la filosofia del "learning-while-doing" ponendo l'enfasi sull'applicazione degli strumenti e della teoria dell'apprendimento a casi di studio reali e di interesse per Bioinformatici e Scienziati dei dati.</i></p> <p><i>Le lezioni frontali fanno affidamento su supporti multimediali (presentazioni PowerPoint) ed hanno come obiettivo l'insegnamento dei fondamenti della teoria dell'apprendimento, delle tecniche di apprendimento supervisionato e delle procedure di validazione dei modelli data-driven.</i></p> <p><i>Le esercitazioni in aula prevedono la scrittura di codice in linguaggio R in grado di eseguire tutte le procedure utili per pulire i dati in ingresso, eseguire analisi preliminari, addestrare e validare algoritmi di apprendimento.</i></p>
<p><b>Risultati di apprendimento previsti</b></p> <p><i>Da indicare per ciascun Descrittore di Dublino (DD=</i></p> <p><b>DD1</b> Conoscenza e capacità di comprensione</p> <p><b>DD2</b> Conoscenza e capacità di comprensione applicate</p> <p><b>DD3-5</b> Competenze trasversali</p>	<p><b>- Descrittore di Dublino 1:</b> <i>conoscenza e capacità di comprensione;</i></p> <ul style="list-style-type: none"> <li>○ Comprensione dei teoremi generali della teoria dell'apprendimento supervisionate e delle loro implicazioni;</li> <li>○ Comprensione delle specificità degli algoritmi di apprendimento automatico introdotte nel corso</li> </ul> <p><b>- Descrittore di Dublino 2:</b> <i>capacità di applicare conoscenza e comprensione;</i></p> <ul style="list-style-type: none"> <li>○ Capacità di applicare le procedure di analisi illustrate durante il corso per condurre analisi su dati omici mai visti prima.</li> </ul> <p><b>- Descrittore di Dublino 3:</b> <i>capacità critiche e di giudizio (occorre indicare le attività che concorrono allo sviluppo di tali abilità.</i></p> <ul style="list-style-type: none"> <li>• <b>Autonomia di giudizio</b></li> </ul> <p><i>Al termine dell'insegnamento lo/la studente/studentessa dovrà essere in grado di</i></p> <ul style="list-style-type: none"> <li>• Capacità di applicare le procedure di analisi illustrate durante il corso per condurre analisi su dati omici mai visti prima.</li> <li>• Capacità di formulare giudizio autonomo circa le performance di classificazione misurate per un certo problema con riferimento a limiti o margini di miglioramento dell'algoritmo di apprendimento proposto</li> <li>• Comprendere implicazioni etiche delle varie fonti di bias nella ricerca accademica.</li> </ul> <p><b>- Descrittore di Dublino 4:</b> <i>capacità di comunicare quanto si è appreso.</i></p> <p>☒ <b>Abilità comunicative</b></p> <p><i>Al termine dell'insegnamento lo/la studente/studentessa dovrà essere in grado di</i></p> <ul style="list-style-type: none"> <li>• Realizzare grafici efficaci nel comunicare esistenza di relazioni associative</li> <li>• Comunicare i risultati scientifici raggiunti ai propri interlocutori con l'ausilio di supporti quali diapositive</li> </ul> <p><b>- Descrittore di Dublino 5:</b> <i>capacità di proseguire lo studio in modo autonomo nel corso della vita. Capacità di apprendere in modo autonomo:</i></p> <p><i>Al termine dell'insegnamento lo/la studente/studentessa dovrà essere in grado di</i></p> <ul style="list-style-type: none"> <li>• Capacità di applicare le procedure di analisi illustrate durante il corso per condurre analisi su dati omici mai visti prima.</li> </ul>

<b>Contenuti di insegnamento (Programma)</b>	<ul style="list-style-type: none"> <li>• <i>Introduction: what is Machine Learning; Problems, data, and tools.</i></li> <li>• <i>An introduction to R.</i></li> <li>• <i>Feature Engineering.</i></li> <li>• <i>Dimensionality reduction; latent space methods; Principal Component Analysis (PCA);</i></li> <li>• <i>Regression problem: linear regression, gradient descent.</i></li> <li>• <i>Classification problems: Overfitting and complexity; bias and variance; training, validation, and test; Cross validation procedure; confusion matrix, accuracy, precision, recall, f1 score.</i></li> <li>• <i>Linear classifiers; Naive Bayes model; Neural Networks, Support Vector Machine.</i></li> <li>• <i>Decision Tree; Ensemble: Bagging, Random forests, Boruta, VSURF.</i></li> <li>• <i>Unsupervised learning: clustering, k-means, hierarchical agglomeration.</i></li> <li>• <i>Hints on Explainable Artificial Intelligence (XAI): SHAP Model.</i></li> <li>• <i>Hints on CNNs and Unet-128.</i></li> <li>• <i>Hints on Cohort Studies.</i></li> </ul>
<b>Testi di riferimento</b>	<b><i>Christopher M. Bishop, Pattern Recognition and Machine Learning</i></b>
<b>Note ai testi di riferimento</b>	
<b>Materiali didattici</b>	<i>Presentazioni PowerPoint su canale Teams.</i>

<b>Valutazione</b>	
<b>Modalità di verifica dell'apprendimento</b>	<p><i>La verifica dell'apprendimento prevede la presentazione di un progetto di analisi di un dataset multi-omico adeguatamente descritto da articolo internazionale revisionato tra pari. Il progetto prevede tre fasi:</i></p> <ul style="list-style-type: none"> <li>• <i>sviluppo software e addestramento di algoritmi</i></li> <li>• <i>misura e validazione delle performance di classificazione/regressione</i></li> <li>• <i>realizzazione di un PowerPoint e colloquio orale di 20 minuti (10 slide) relativamente alle performance della catena di analisi proposta</i></li> </ul> <p><i>Il colloquio orale ha la finalità di comprendere se lo studente ha acquisito la capacità di applicare autonomamente gli strumenti descritti e dimostrati in aula a problemi scientifici aperti.</i></p> <p><i>Tali competenze verranno comprovate durante il colloquio finale con l'ausilio di una presentazione PowerPoint. La presentazione PowerPoint prevede alcune slide introduttive sul problema scientifico oggetto di indagine, sulla catena di analisi e sui metodi impiegati, ed infine, alcune slide conclusive relativamente a risultati ottenuti e limiti delle procedure adottate rispetto al particolare dato investigato.</i></p> <p><i>In particolare, la modalità di valutazione dell'esame farà particolare attenzione a misurare la capacità dello studente di individuare relazioni di associazione statistica in dati di elevata dimensionalità.</i></p>

<p>Criteri di valutazione</p>	<ul style="list-style-type: none"> <li>• <i>Conoscenza e capacità di comprensione:</i> <ul style="list-style-type: none"> <li>• <i>Comprensione dei teoremi generali della teoria dell'apprendimento supervisionate;</i></li> </ul> </li> <li>• <i>Conoscenza e capacità di comprensione applicate:</i> <ul style="list-style-type: none"> <li>• <i>Comprensione implicazione dei teoremi generali della teoria dell'apprendimento supervisionato (VC-dimension, bias-variance trade-off)</i></li> <li>• <i>Comprensione specificità dei vari algoritmi di apprendimento</i></li> </ul> </li> <li>• <i>Autonomia di giudizio:</i> <ul style="list-style-type: none"> <li>• <i>Capacità di valutare la qualità della propria ricerca da un punto di vista statistico e della replicabilità dei risultati</i></li> </ul> </li> <li>• <i>Abilità comunicative:</i> <ul style="list-style-type: none"> <li>• <i>Capacità di comunicare i risultati scientifici della propria ricerca</i></li> </ul> </li> <li>• <i>Capacità di apprendere:</i> <ul style="list-style-type: none"> <li>• <i>Capacità di applicare gli strumenti illustrati a dati nuovi in completa autonomia.</i></li> </ul> </li> </ul>
<p>Criteri di misurazione dell'apprendimento e di attribuzione del voto finale</p>	<p><i>Il voto finale è attribuito in trentesimi. L'esame prevede una sola valutazione orale e si intende superato quando il voto è maggiore o uguale a 18.</i></p>

General information	
Year of the course	2023/2024
Academic calendar (starting and ending date)	march 2024 –may 2024
Credits (CFU/ETCS):	6
SSD	Fis/07 o ING-INF/05
Language	Italian or English
Mode of attendance	Optional

Professor/ Lecturer	
Name and Surname	Roberto Cilli
E-mail	<a href="mailto:roberto.cilli@uniba.it">roberto.cilli@uniba.it</a>
Telephone	327 730 8761
Department and address	Physics Department, Room 124
Virtual room	Italian or English
Office Hours (and modalities: e.g., by appointment, on line, etc.)	Every Friday, from 15 to 19 – by appointment

Work schedule			
Hours			
Total	Lectures	Hands-on (laboratory, workshops, working groups, seminars, field trips)	Out-of-class study hours/ Self-study hours
150	32	30	88
CFU/ETCS			
6	4	2	

<b>Learning Objectives</b>	<i>The course will provide students with the necessary skills to understand and implement machine learning algorithms for knowledge extraction in large multi-omics data sets.</i>
<b>Course prerequisites</b>	<i>There are no specific prerequisites that differs from those required for admission to the master's degree course. However, prior knowledge of R-language is considered as an asset.</i>

<b>Teaching strategies</b>	<p><i>The teaching strategy follows the "learning-while-doing" philosophy with emphasis on the application of tools and learning theory to real-world case studies of interest to Bioinformaticians and Data Scientists. Therefore, the lessons consists in face-to-face lectures alternated coding classes.</i></p> <p><i>Classroom lectures rely on multimedia aids (PowerPoint presentations) and aim to teach the fundamentals of learning theory, supervised learning techniques, and data-driven model validation procedures.</i></p> <p><i>Coding classes are devoted to code R scripts performing all useful procedures for cleaning input data, performing preliminary analysis, and training and validating learning algorithms.</i></p>
----------------------------	--

<b>Expected learning outcomes in terms of</b>	
<b>Knowledge and understanding on in terms of the Dublin descriptors</b>	<p><b>Dublin Descriptor 1:</b> Knowledge and understanding skills;</p> <ul style="list-style-type: none"> <li>• Understanding of the general theorems of supervised learning theory and their implications.</li> <li>• Understanding of the specifics of machine learning algorithms introduced in the course.</li> </ul> <p><b>Dublin Descriptor 2:</b> Ability to apply knowledge and understanding.</p> <ul style="list-style-type: none"> <li>• Ability to apply the analysis procedures explained in the course to conduct analyses on omics data never seen before.</li> </ul> <p><b>Dublin Descriptor 3:</b> Critical and judgment skills. After completing, the student is expected to be able to:</p> <ul style="list-style-type: none"> <li>• apply the analysis procedures outlined in the course to conduct analyses on omics data never seen before.</li> <li>• make independent judgment about the classification performance measured for a given problem with reference to limitations or room for improvement of the proposed learning algorithm.</li> <li>• Understand ethical implications of various sources of bias in academic research.</li> </ul> <p><b>Dublin Descriptor 4:</b> Ability to communicate what has been learned. After completing, the student is expected to have learnt to:</p> <ul style="list-style-type: none"> <li>• Communicate scientific results to their peers or stakeholders by using effective graphs.</li> </ul> <p><b>Dublin Descriptor 5:</b> Ability to learn independently. After completing, the student is expected to have learnt to:</p> <ul style="list-style-type: none"> <li>• How to conduct empirical analyses on omics data never seen before by following the procedure outlined during the course</li> </ul>
<b>Syllabus</b>	
<b>Content knowledge</b>	<ul style="list-style-type: none"> <li>• <i>Introduction: what is Machine Learning; Problems, data, and tools.</i></li> <li>• <i>An introduction to R.</i></li> <li>• <i>Feature Engineering.</i></li> <li>• <i>Dimensionality reduction; latent space methods; Principal Component Analysis (PCA);</i></li> <li>• <i>Regression problem: linear regression, gradient descent.</i></li> <li>• <i>Classification problems: Overfitting and complexity; bias and variance; training, validation, and test; Cross validation procedure; confusion matrix, accuracy, precision, recall, f1 score.</i></li> <li>• <i>Linear classifiers; Naive Bayes model; Neural Networks, Support Vector Machine.</i></li> <li>• <i>Decision Tree; Ensemble: Bagging, Random forests, Boruta, VSURF.</i></li> <li>• <i>Unsupervised learning: clustering, k-means, hierarchical agglomeration.</i></li> <li>• <i>Hints on Explainable Artificial Intelligence (XAI): SHAP Model.</i></li> <li>• <i>Hints on CNNs and Unet-128.</i></li> <li>• <i>Hints on Cohort Studies.</i></li> </ul>
<b>Texts and readings</b>	<b>Christopher M. Bishop, <i>Pattern Recognition and Machine Learning</i></b>
<b>Notes, additional materials</b>	<i>PowerPoint presentations</i>
<b>Repository</b>	

Assessment	
Assessment methods	<p><i>Learning will be assessed by involving the students in a personal project consisting in conducting a data-driven analysis of omics data and by presenting the obtained results. The omics data should be adequately described by an international peer-reviewed article. More precisely, the project involves three steps:</i></p> <ul style="list-style-type: none"> <li>- software development and algorithm training</li> <li>- measurement and validation of classification/regression performance</li> <li>- creation of a 20-minute PowerPoint and oral interview (10 slides) regarding the performance of the proposed analysis chain</li> </ul> <p><i>The oral interview aims to understand whether the student has acquired the ability to independently apply the AI/ML tools and to defend their results. In particular, the assessment mode of the exam will pay special attention to measuring the student's ability to identify statistical association relationships in data of high dimensionality.</i></p>
Assessment criteria	<p><i>The student will be evaluated according to the following pillars.</i></p> <p><i>Knowledge and understanding skills:</i></p> <ul style="list-style-type: none"> <li>• <i>Understanding of the general theorems of supervised learning theory;</i></li> </ul> <p><i>Applied knowledge and understanding skills:</i></p> <ul style="list-style-type: none"> <li>• <i>Understanding implication of general theorems of supervised learning theory (VC-dimension, bias-variance trade-off)</i></li> <li>• <i>Understanding specificity of various learning algorithms</i></li> </ul> <p><i>Autonomy of judgment:</i></p> <ul style="list-style-type: none"> <li>• <i>Ability to assess the quality of one's research from a statistical perspective and replicability of results</i></li> </ul> <p><i>Communication skills:</i></p> <ul style="list-style-type: none"> <li>• <i>- Ability to communicate scientific results to their peers</i></li> </ul> <p><i>Ability to learn:</i></p> <ul style="list-style-type: none"> <li>• <i>- Ability to apply the tools explained to new data with complete autonomy.</i></li> </ul>
Final exam and grading criteria	<p><i>The examination involves only one oral assessment and is considered passed when the grade is greater than 18/30.</i></p>