

<b>CORSO DI STUDIO</b>	<b>LAUREA IN MATEMATICA (L-35)</b>
<b>ANNO ACCADEMICO</b>	<b>2023-2024</b>
<b>INSEGNAMENTO</b>	<b>METODI NUMERICI IN DATA SCIENCE</b>

Principali informazioni sull'insegnamento	
Anno di corso	Terzo
Periodo di erogazione	Secondo semestre (26 febbraio 2024 – 31 maggio 2024)
Crediti formativi universitari (CFU)	7
Settore scientifico disciplinare (SSD)	MAT/08 – Analisi Numerica
Lingua di erogazione	Italiano
Modalità di frequenza	Facoltativa

Docenti		
Nome e cognome	Flavia Esposito (titolare)	Nicoletta Del Buono
Indirizzo mail	flavia.esposito@uniba.it	nicoletta.delbuono@uniba.it
Telefono	+39 080 544 2711	+39 080 544 2711
Sede	Dipartimento di Matematica stanza 24 secondo piano	Dipartimento di Matematica stanza 24 secondo piano
Sede virtuale	687hs39	687hs39
Pagina web	<a href="https://www.dm.uniba.it/it/members/esposito">https://www.dm.uniba.it/it/members/esposito</a>	<a href="https://www.dm.uniba.it/it/members/delbuono">https://www.dm.uniba.it/it/members/delbuono</a>
Ricevimento	Su appuntamento, da concordare per e-mail	Su appuntamento, da concordare per e-mail

Organizzazione della didattica				
	Totali	Didattica frontale	Pratica (esercitazioni/laboratori)	Studio individuale
<b>Ore</b>	175	48	15	112
<b>CFU</b>	7	6	1	

Obiettivi formativi	
	Acquisizione delle tecniche numeriche di base per le applicazioni di data science. Acquisizione delle tecniche numeriche di base per la Exploratory Data Analysis e dei metodi di algebra lineare numerica per il trattamento di alcuni problemi di apprendimento dai dati.

Prerequisiti	
	Le conoscenze acquisite nella laurea della classe L-35 con riferimento particolare alle discipline di Calcolo Numerico e della Analisi Matematica classica in una e più variabili.

Syllabus	
Contenuti dell'insegnamento (Programma)	-Introduzione ai metodi di fattorizzazione per la Data Science. Algebra Lineare e apprendimento da dati strutturati. Rappresentazione di dati attraverso vettori e matrici. Decomposizione a valori singolari di una matrice di dati. Proprietà dei vettori singolari destri e sinistri e loro relazione con gli spazi vettoriali fondamentali generati da una matrice di dati. L'importanza della SVD nella Data Science: esempi. La trasformazione di Karhunen-Loève e la sua relazione con la SVD. Componenti principali e



migliore approssimazione low-rank di una matrice di dati. Forma ridotta della SVD e teorema di Eckart-Young. Analisi di dati attraverso la fattorizzazione SVD della matrice di covarianza e di correlazione. Analisi delle componenti principali e equivalenza con il problema di massimo per la varianza dei dati. La geometria della PCA. Analisi di dati nonnegativi. Matrici positive e nonnegative. Teoremi di Perron-Frobenius. Modello Eigenface.

-Il modello vettoriale dell'informazione (VSM) e il Latent Semantic Indexing. Introduzione alla formalizzazione algebrica dei problemi di information retrieval. Processo di indexing automatico, operazioni di stop-listing e stemming. Funzioni di pesatura degli index-term. Applicazione delle fattorizzazioni QR e SVD alla matrice termini-documenti e loro interpretazione geometrica. Approssimazione low-rank dello spazio semantico. Il processo di query-matching e misure di similarità. Confronto termine-termine e clustering di termini sinonimi. Accenno ai meccanismi di relevance feedback basati sull'utilizzo della SVD troncata della matrice termini-documenti.

-Modelli basati su autovalori e autovettori per web information retrieval. La struttura a hyperlink del web e sua rappresentazione attraverso i grafi non orientati. Concetti di inlink e outlink. Modelli HITS e PageRank per il ranking di reti web. Costruzione delle matrici di Hub e Authority e loro proprietà. Costruzione della matrice di Google e sue proprietà. Riducibilità e grafi. Convergenza del modello PageRank e irriducibilità della matrice di Google.

-Introduzione a problemi su reti e grafi. Esempi e problemi (problema dei ponti di Königsberg e delle porte di una casa). Definizioni di cammini e cicli. Cammini euleriani e hamiltoniani. Grafi e rappresentazioni matriciali. Matrici di adiacenza, incidenza, matrice dei gradi e matrice laplaciana  $L_G$  di un grafo. Alcune proprietà spettrali della matrice di adiacenza e della matrice  $L_G$ . Autovalori e misure di connettività di un grafo. Cenni su grafi completi e loro proprietà. Problema del cammino di costo minimo e algoritmo di Dijkstra.

-Introduzione al Machine Learning (problemi supervisionati e non supervisionati, introduzione alla classificazione e alla regressione con esempi, concetti di over e under fitting, trade-off bias-varianza), Struttura del dato (vettori, matrici e tensori), analisi esplorativa del dato (EDA), tipi di variabile e cenni di statistica descrittiva (misure di tendenza centrale, misure di variabilità, quantili) con relative rappresentazioni grafiche (box-plot, diagramma a barre, istogrammi, scatter plot). Ispezionare le relazioni tra variabili, Definizione del coefficiente di correlazione di Pearson.

-Pre-processing del dato: definizione/classificazione e trattamento di missing values, definizione outliers e identificazione outliers, trasformazioni e normalizzazioni del dato. Il problema della Curse of Dimensionality e come risolverlo con tecniche di fattorizzazione lineare, Interpretazioni delle decomposizioni matriciali, Decomposizioni matriciali come sistemi di raccomandazioni (cenni).

-Introduzione alla Nonnegative Matrix Factorization (NMF). Storia ed esempi. Motivazioni e interpretazione dei fattori non negativi, rappresentazione part-based e learning. Esempio eigenfaces e confronto tra VQ, PCA e NMF, Formalizzazione della NMF come problema di ottimizzazione matriciale. Analisi delle due funzioni obiettivo più utilizzate, interpretazione probabilistica della NMF. Divergenze Beta e di Bregman (definizioni ed equivalenze). Interpretazione geometrica della NMF. Nonnegative Rank Factorization. Algoritmi numerici per il calcolo della NMF

	<p>(Multiplicative Updates e Block Coordinate Descent). Teoremi di convergenza e dimostrazioni. Cenni sulla NMF regolarizzata.</p> <p>-Introduzione al problema di clustering, distanze utili per il clustering, Clustering Gerarchico divisivo e agglomerativo (rappresentazioni tramite dendrogramma e metodi di linkaggio: singolo, completo, medio), K-means. Metodi euristici per la scelta del numero ottimale di clusters nel K-means, Equivalenza tra NMF e K-means con vincolo rilassato di ortogonalità (Teoremi e dimostrazioni), K-medioide (cenni). Indici Interni, Indici Esterni e altre misure di bontà degli algoritmi di clustering. Algoritmi Model e Density Based (Mixture Model e DBSCAN): cenni.</p> <p>Introduzione ai problemi supervisionati. Rischio Atteso e Rischio Atteso Empirico. Introduzione alla regressione, problema di under e overfitting nella regressione. Regressione Lineare Semplice e Multipla, Stima dei coefficienti con OLS (Ordinary Least Squares ed equazioni Normali). Teorema di Gauss-Markov e condizioni per il miglior stimatore lineare corretto. Interpretazione geometrica della regressione. Regressione per lo studio della collinearità e l'individuazione degli outliers in matrici di dati. Misure di valutazione dei modelli di regressione (errori SST, SSE, SSR, <math>R^2</math>, <math>R^2</math> adjusted, RMSE, MAE, F statistica). Regressione polinomiale (cenni).</p> <p>-Introduzione ai problemi di classificazione. SVD per classificazione handwritten digits, K-Nearest Neighbors, Approcci ad albero per la classificazione (Alberi Decisionali e Random Forest). Support Vector Machines, Maximal Margin Classifier, Soft Margin Classifier, Teorema di Mercer, Misure per la bontà della Classificazione, Cross Validation (cenni).</p> <p>-Esercitazione in R degli argomenti trattati</p>
Testi di riferimento	<p>-G. Strang, Linear Algebra and Learning from Data, Wellesley-Cambridge Press, 2019</p> <p>- C. Meyer, Matrix Analysis and Applied Linear Algebra, SIAM, 2003.</p> <p>- I.T. Jolliffe, Principal Component Analysis, Second Edition, Springer, 2002</p> <p>- A. Cichocki, R. Zdunek, A.H. Phan, S.I Amari, Nonnegative Matrix and Tensor Factorizations, Wiley, 2009</p> <p>- T. Hastie, R. Tibshirani J. Friedman: The Elements of Statistical Learning: Data Mining, Inference, and Prediction. Second Edition, 2009</p> <p>- N. Gillis. Nonnegative Matrix Factorization, SIAM, 2020</p>
Note ai testi di riferimento	
Materiali didattici	Appunti e riferimenti specifici forniti dal docente

<b>Risultati di apprendimento previsti (secondo i Descrittori di Dublino)</b>	
DD1 Conoscenza e capacità di comprensione	<p>Acquisizione di alcune tecniche per la risoluzione di problemi di classificazione e clustering.</p> <p>Capacità di utilizzare codici numerici che implementano tecniche standard per l'analisi di dati reali.</p>
DD2 Conoscenza e capacità di comprensione applicate	<p>Le conoscenze teoriche e pratiche acquisite vengono utilizzate in campi matematici applicati e per risolvere problemi nel contesto dell'apprendimento dai dati.</p>
DD3-5 Competenze trasversali	<p><i>DD3 Autonomia di giudizio:</i> Capacità di individuare le tecniche numeriche adatte per affrontare e risolvere numericamente i problemi derivanti da applicazioni reali che coinvolgono i big data.</p>
	<p><i>DD4 Abilità comunicative:</i> Acquisizione del linguaggio e del formalismo matematico necessario per la consultazione e comprensione dei testi, l'esposizione delle conoscenze acquisite, la descrizione, l'analisi e la risoluzione dei problemi applicativi e di Data Science</p>

*DD5 Capacità di apprendere:* Acquisizione di un metodo di studio adeguato, supportato dalla consultazione dei testi e dalla implementazione al calcolatore delle tecniche numeriche esposte durante il corso.

Metodi didattici	
	- Lezioni frontali condotte con l'ausilio di supporti didattici (slide). - Esercitazioni al calcolatore.

Valutazione	
Modalità di verifica dell'apprendimento	Prova orale sul programma svolto nel corso delle lezioni ed esercitazioni o progetto assegnato dal docente
Criteri di valutazione	<ul style="list-style-type: none"> <li>• <i>Conoscenza e capacità di comprensione:</i> Gli studenti devono dimostrare una adeguata conoscenza dei contenuti dell'insegnamento</li> <li>• <i>Conoscenza e capacità di comprensione applicate:</i> Gli studenti devono dimostrare una adeguata conoscenza delle possibili applicazioni dei concetti teorici e possedere una adeguata capacità di implementare tali applicazioni</li> <li>• <i>Autonomia di giudizio:</i> Gli studenti devono dimostrare una adeguata autonomia nella selezione dei concetti teorici più idonei alla risoluzione di problemi pratici</li> <li>• <i>Abilità comunicative:</i> Gli studenti devono dimostrare una adeguata capacità espositiva dei contenuti studiati e una adeguata capacità di analisi e sintesi</li> <li>• <i>Capacità di apprendere:</i> Gli studenti devono dimostrare una buona capacità di effettuare collegamenti interdisciplinari</li> </ul>
Criteri di misurazione dell'apprendimento e di attribuzione del voto finale	<p>Il voto finale è attribuito in trentesimi. L'esame si intende superato quando il voto è maggiore o uguale a 18. Per la formulazione del voto finale si prenderanno in considerazione i seguenti indicatori: grado di conoscenza dei contenuti e degli argomenti dell'insegnamento, capacità e correttezza nell'applicare i concetti fondamentali trattati durante le lezioni frontali e le esercitazioni, qualità della esposizione orale.</p> <p>Tutti gli argomenti del programma contribuiscono in modo uguale alla formulazione del voto finale</p>

Ulteriori informazioni	
	La frequenza delle lezioni ed esercitazioni è fortemente consigliata