

COURSE OF STUDY	TWO-YEAR MASTER OF SCIENCE PROGRAMME IN MATHEMATICS
ACADEMIC YEAR	2023-2024
ACADEMIC SUBJECT	STATISTICS FOR MACHINE LEARNING

General information	
Term	Second semester (February 26, 2024 – May 31, 2024)
European Credit Transfer and Accumulation System credits (ECTS)	7
SSD	MAT/06 – Probability and Mathematical Statistics
Language	Italian
Mode of attendance	Not mandatory

Lecturer	
Name and surname	Marcello De Giosa
E-mail	marcello.degiosa@uniba.it
Telephone	+39 080 544 2707
Department and office	Department of Mathematics, floor 4, room 12.
Virtual meeting room	
Web page	https://www.dm.uniba.it/it/members/degiosa
Office hours	Wednesday, from 10 a.m. to 12 a.m., by email appointment.

Work schedule				
	Total	Lectures	Hands-on learning (recitations)	Self-study
Hours	175	48	15	112
ECTS credits	7	6	1	

Learning objectives	
	<p>Acquiring knowledge of the main statistical predictive methods of Statistical Machine Learning, and their correct computational implementation with the R programming language e the package collection tidymodels.</p> <p>Particular attention is paid to the theoretical contributions of Statistics and to good methodological practices to produce high quality statistical models of Machine Learning and to correctly interpret and present the conclusions of the analysis of complex data structures.</p>

Course prerequisites	
	Calculus, multivariate calculus, linear algebra, elementary probability, as usually offered during a degree in Mathematics of L-35 class.

Syllabus	
Course contents	<p>Introduction to Statistical Machine Learning. Trade-off between prediction accuracy and model interpretability. Supervised and unsupervised learning. Regression versus classification. Trade-off between bias and variance. Introductory R lab.</p> <p>Linear regression (LR). Simple and multiple LR. Coefficients estimation and accuracy. Qualitative predictors. Comparison with KNN. Interactions. R lab on</p>

	<p>LR.</p> <p>Classification. Logistic Regression. The simple and multiple cases. Generative classification methods: linear and quadratic discriminant analysis. Naïve Bayes. Comparing classification methods. Generalized linear models: Poisson regression. R lab on classification.</p> <p>Resampling methods. Cross-validation and bias-variance tradeoff. Cross validation in classification. The bootstrap. R lab on resampling.</p> <p>Linear model regularization. Shrinkage: Ridge and Lasso. Hyperparameters tuning. Dimension reduction: principal component regression and partial least squares. High dimensional data. R lab on regularization.</p> <p>Beyond linearity. Polynomial regression, step functions, basis functions. Splines regression and smoothing splines. Generalize additive models for regression and classification. R lab on nonlinearity.</p> <p>Decision tree. Regression and classification regression trees. Bagging, Random Forests, Boosting, Bayesian additive regression trees (BART). R lab on decision tree methods.</p> <p>Support vector machines. Maximal margin classifier. Support vectors classifier. Support vector machines. More than two classes: one-versus-one and one-versus-all methods. Comparing with logistic regression. R lab on support vector machines.</p> <p>Deep Learning. Single Layer Neural Networks. Multilayer Neural Networks. Convolutional Neural Networks: Convolution layer, pooling layer, architecture, data augmentation, pretrained classifier. Backpropagation. Regularization and Stochastic Gradient Descent. Dropout. Tuning. Interpolation and Double Descent. R lab on deep learning with Keras and Torch.</p> <p>Unsupervised Learning. Principal Component Analysis (PCA). Clustering: K-means, hierarchical clustering. R lab on unsupervised learning.</p> <p>Introduction to tidymodels. Review of R lab with tidymodels.</p> <p>Review of resampling methods. Resubstitution. Cross-validation. Repeated Cross-validation. Monte-Carlo Cross-validation. Bootstrap. Performance estimation. Parallel processing.</p> <p>Model and parameters tuning. Grid search and iterative search. Bayesian optimization. Simulated annealing.</p> <p>Review of dimension reduction methods. Principal Components Analysis (PCA), Partial Least Squares (PLS), Independent Component Analysis (ICA), Uniform Manifolds Approximation and Projection (UMAP). Comparison of models based on different dimension reduction methods. R lab on dimensions reduction.</p>
Reference books	<p>G. James, D. Witten, T. Hastie, R. Tibshirani – An introduction to Statistical Learning with application in R – 2023 – Springer</p> <p>T. Hastie, R. Tibshirani, J. Friedman – The Elements of Statistical Learning. Data Mining, Inference, and Prediction – 2009 – Springer</p> <p>K.P. Murphy – Probabilistic Machine Learning – 2022 – MIT Press</p> <p>J. Fan, R. Li, C. Zhang, H. Zou – Statistical Foundations of Data Science – 2020 – CRC Press</p> <p>M. Kuhn, k. Johnson – Applied Predictive Modeling – 2016 – Springer</p> <p>D.P. Kroese, Z.I. Botev, T. Taimre, R. Vaisman - Data Science and Machine Learning. Mathematical and Statistical Methods. - 2020 - CRC Press</p> <p>R. Irizarry – Introduction to Data Science – 2019 – CRC Press</p>
Additional course materials	
Repository	

Expected learning outcomes	
Knowledge and understanding	Acquiring the main concepts and techniques, and R code understanding for Statistical Machine Learning.
Applying knowledge and understanding	Being able to apply the theoretic results and R code to solve real world Statistical Machine Learning problems.
Soft skills	<i>Making judgements</i> : Being able to choose the right Statistical Machine Learning tools and R code.
	<i>Communication skills</i> : Being able to report and clearly explain the Statistical Machine Learning methods and results.
	<i>Learning skills</i> : Being able to apply learning methods to improve knowledge in the field of Statistical Machine Learning.

Teaching methods	
	In person frontal theoretical lessons and R laboratory with use of slides.

Assessment	
Assessment methods	The final exam consists of an oral part and an R laboratory part.
Evaluation criteria	<ul style="list-style-type: none"> • <i>Knowledge and understanding</i>: mastering and deep understanding of the Statistical Machine Learning methods and their translation in modern R code. • <i>Applying knowledge and understanding</i>: being able to apply the Statistical Machine Learning methods and related R code to solve practical real world problems. • <i>Making judgement</i>: mastering the ability to choose the right Statistical Machine Learning methods and R code to solve real problems. • <i>Communication skills</i>: ability to explain the applied methods and their results in a clear way. • <i>Learning skills</i>: well organizing the acquired knowledge and being able to autonomously find sources of insight.
Grading policy	The exam consists in the discussion of some methods among those proposed during the classes and the discussion of some R code chosen at random from those presented during the classes. The vote, out of thirty, will depend on the mastery of the topics, knowledge of the codes, but also on the ability to clearly explain results.

Further information	



UNIVERSITÀ
DEGLI STUDI DI BARI
ALDO MORO

CONSIGLIO INTERCLASSE
IN MATEMATICA