

## Exploratory Data Analysis and Modeling

### 1. Course/Module information

<b>Faculty</b>	CIM				
<b>Department</b>	Software Engineering and Automation				
<b>Study cycle</b>	Master's Degree				
<b>Study program</b>	Data Science				
<b>Year of study</b>	<b>Semester</b>	<b>Evaluation type</b>	<b>Formative category</b>	<b>Optionality category</b>	<b>ECTS credits</b>
1 <sup>st</sup> Year ( <i>full-time education</i> )	1	E	F	O	5

### 2. Estimated total time

Total hours in the curriculum plan	Including				
	Auditory hours		Individual work		
	Lecture	Laboratory/ seminar	Term paper	Study of theoretical material	Application development
Full-time education	20	20	-	110	-

### 3. Prerequisites for access to the course/module

<b>According to the curriculum plan</b>	Mathematics (Linear Algebra, Differential and Integral Calculus, Probability Theory, Mathematical Statistics), Data Structures and Algorithms, Programming (Python or R), Databases, Web Technologies.
<b>According to competencies</b>	Basic data-related skills and knowledge, fundamental statistical concepts, key concepts from higher mathematics (including linear algebra, differential and integral calculus, etc.), data manipulation skills, understanding graphical representations, critical thinking and logical reasoning, basic programming skills (programming concepts in any language, preferably Python or R), data collection concepts, knowledge of how data is stored and extracted from databases, interest and curiosity for data analysis, understanding of the business context, self-learning skills.

### 4. Conditions for conducting the educational process

<b>Lecture</b>	A projector and computer are required to present theoretical material in the classroom. Student tardiness, as well as phone conversations during the lecture, will not be tolerated.
<b>Laboratory/ seminar</b>	Students will engage in seminars under the guidance of the professor and assistant and will complete reports according to the methodological instructions. The deadline for submitting individual work is one week from completion. If the work is submitted late, a penalty of 1 point per week of delay will be applied.

### 5. Specific competencies acquired

<b>Professional competencies</b>	<b>CPM1</b> System architecture design and development. <b>CPM2</b> Monitoring technological trends. Innovation. Sustainable development. <b>CPM3</b> Application development. Component integration. Systems engineering.
<b>Transversal competencies</b>	<b>CT1.</b> Autonomy and responsibility <b>CT2.</b> Social interaction <b>CT3.</b> Personal and professional development

## 6. Course/Module objectives

<b>General objective</b>	The development of knowledge and the formation of skills necessary for the effective use of data analysis methods and techniques to detect and capture the relationships and patterns hidden behind the data, and to obtain a simplified, clear, and easily interpretable representation of these relationships and patterns.
<b>Specific objectives</b>	<p>To achieve the general objective, it is necessary to develop knowledge and skills in:</p> <ul style="list-style-type: none"> <li>• Data Profiling and Summarization: Techniques for effective profiling and summarizing of datasets.</li> <li>• Data Visualization: Skills for visually representing data using various charts, graphs, etc.</li> <li>• Pattern and Trend Identification: Identifying patterns, trends, and relationships within data to generate insights.</li> <li>• Handling Missing Data and Outliers: Methods for handling missing data and detecting outliers during exploratory analysis.</li> <li>• Statistical Analysis Techniques: Statistical methods to analyze data and draw conclusions.</li> <li>• Preprocessing Data: Knowledge of data preprocessing techniques to prepare data for analysis.</li> <li>• Domain-Specific EDA: EDA techniques specific to various fields, such as finance, healthcare, or marketing.</li> <li>• Communicating Findings: Skills to effectively communicate insights and findings from EDA to both technical and non-technical audiences.</li> <li>• Collaboration Practices in EDA: Collaborating with colleagues in performing EDA and sharing information for collective understanding.</li> <li>• Ethical and Responsible EDA: Ethical considerations and responsibilities related to performing EDA, including data privacy and protection.</li> </ul>

## 7. Course/Module content

Syllabus of teaching activities	Number of hours
<b>Course topics</b>	
<b>T1 EDA Research Topic, motivation and methods</b> Descriptive statistics and graphical techniques for data exploration. Definitions of basic concepts such as population, sample, and observation. Measures of central tendency (mean, mode, median) and their properties. Measures of variability and dispersion (range, variance, standard deviation, skewness, kurtosis) and their properties. Chebyshev's theorem and the empirical rule of distribution for normal distribution	2
<b>T2 Types of Data and graphical methods of exploration.</b> Qualitative data: Nominal vs. ordinal data. Quantitative data: Discrete vs. continuous data. Exploring qualitative data - frequency charts and pie charts. Exploring quantitative data - histograms, quartiles, percentiles, and measures of relative position. Interquartile Range and Box Plots. Time series data - line charts and trends.	2
<b>T3. Relationships between two variables.</b> Direction, Type, and Relevance of the Relationship. Correlation, Covariance and Pearson's Correlation Coefficient. Interpretation and Examples. False Correlation.	2
<b>T4. Least squares method.</b> Simple linear regression. Sum of squared residuals and stochastic error. Residual minimization problem. Examples.	2
<b>T5. Multiple linear regression and fit quality.</b> R-squared, total sum of squares, residual sum of squares and explained sum of squares.	2
<b>T6. Hypothesis testing.</b> Null and alternative hypotheses. Type I and Type II errors. Decision rule. Acceptance and rejection regions. T-statistic, critical t-value, and p-value for significance level. Examples.	2
<b>T7. Multiple hypothesis testing and regression significance.</b> F-test, critical F-value, degrees of freedom.	2

Syllabus of teaching activities	Number of hours
<b>T8. Model specification and predictor selection.</b> OLS assumptions. Influence of omitted variables. Suppressing the intercept. Polynomial regressions. Dummy variables.	2
<b>T9. Logistic regression.</b> Understanding logistic regression and the Logit function. Interpreting logistic regression results.	2
<b>T10. Time series.</b> Trend, cyclical, seasonality. AR, MA, ARIMA, SARIMA, VARIMA. Difference and stationarity. Dickey-Fuller (DF) test and Augmented Dickey-Fuller (ADF) test. Autocorrelation function (ACF) and partial autocorrelation function (PACF).	2
<b>Total course, hours</b>	<b>20</b>
<b>Practical work topics</b>	
<b>Practical work No. 1.</b> Introduction to Using Python in Google Collaboratory - Loading datasets, calculating, and exploring descriptive statistics.	2
<b>Practical work No. 2.</b> Data visualization and analysis using graphical methods in Python.	2
<b>Practical work No. 3.</b> Correlations between variables in Python. Heatmaps and their interpretation.	2
<b>Practical work No. 4.</b> Implementation and interpretation of a simple linear regression model in Python.	2
<b>Practical work No. 5.</b> Multivariate regression in Python and evaluating the goodness of fit of the model.	2
<b>Practical work No. 6.</b> Hypothesis testing in Python, interpreting statistical significance.	2
<b>Practical work No. 7.</b> Multivariate regression in Python and joint significance testing using the F-test.	2
<b>Practical work No. 8.</b> Model specification testing and simulating omitted variables.	2
<b>Practical work No. 9.</b> Logistic regression in Python. Encoding binary variables.	2
<b>Practical work No. 10.</b> Time series modeling in Python: AR, MA, ARIMA. ACF and PACF. DF and ADF tests in Python – implementation and interpretation.	2
<b>Total practical work, hours</b>	<b>20</b>

## 8. Using generative AI

<b>Permission to use</b>	<p>The use of generative AI in assignments and projects is permitted, provided that students adhere to the following rules:</p> <ul style="list-style-type: none"> <li>• Generative AI may be used to generate ideas, text structures, or code, but all generated materials must be reviewed and adjusted by the student to ensure that they meet academic requirements.</li> <li>• Any use of generative AI must be declared in the appendix section of each paper, using the phrase: "During the preparation of this paper, the author used [NAME OF TOOL / SERVICE] for the purpose of [REASON]. After using this tool / service, the author reviewed and edited the content as necessary and assumes full responsibility for the content of the paper."</li> </ul>
<b>Restrictions to use</b>	<p>Students <b><i>MUSTN'T consider generative AI as a reliable source of information</i></b>, as it does not provide clear references or documented sources.</p> <ul style="list-style-type: none"> <li>• <b><i>Direct citation of AI-generated content</i></b> in academic papers as if it were a primary source <b><i>isn't permitted</i></b>.</li> <li>• Activities in which the use of <b>generative AI is prohibited</b> are specified by the teacher and are usually <b><i>intermediate and final assessments</i></b> or that don't involve professional competence development activities.</li> </ul>

## 9. Bibliographic references

<b>Main</b>	<ol style="list-style-type: none"> <li>1. Keller Gerald, Statistics for Management and Economics, South-Western College Publishing; 10th edition, 2014</li> <li>2. Studenmund Arnold H., Using Econometrics: A practical Guide, Pearson Education, 2014</li> <li>3. Educational Materials and Bibliographic Sources on FCIM's ELSE Platform: <a href="https://else.fcim.utm.md/enrol/index.php?id=701">https://else.fcim.utm.md/enrol/index.php?id=701</a></li> </ol>
<b>Supplementary</b>	<ol style="list-style-type: none"> <li>1. Wes McKinney, Python for Data Analysis, O'Reilly Media, 3rd Edition, 2022 ISBN: 9781098104009</li> <li>2. Jake Vanderplas, Python Data Science Handbook: Essential Tools for Working with Data. O'Reilly Media, 2nd Edition, 2023, ISBN-10: 1098121228, ISBN-13: 978-1098121228</li> </ol>

## 10. Evaluation

Periodic		Current	Individual study	Project/thesis	Exam
PE 1	PE 2				
Full-time education					
15%	15%	15%	15%	-	40%
Minimum performance standards Attendance at lectures; activity and quality of preparation for lectures and practical works; Obtaining a minimum grade of "5" for each assessment and practical work; Demonstrating knowledge of the theoretical content of the course and the Python/R language in the final exam paper.					

## 11. Evaluation criteria

Activity	Evaluation components	Evaluation method, Evaluation criteria	Weight in the final grade of the activity	Weight in the course evaluation
<b>Full-time education</b>				
<b>Periodic evaluation I</b>	Theoretical content, topics 1-4	Test on MOODLE	100%	<b>15%</b>
<b>Periodic evaluation II</b>	Theoretical content, topics 6-8	Test on MOODLE	100%	<b>15%</b>
<b>Current evaluation</b>	Practical activity	Discussions during seminars	50%	<b>15%</b>
		File completed with Reports for each Case Study under discussion	50%	
<b>Individual study</b>	Research on the topic	Presentation/public speech	100%	<b>15%</b>
<b>Final evaluation</b>	Theoretical and practical content	Oral exam. Grading according to grading scale	100%	<b>40%</b>