



Universitat de Lleida

# DEGREE CURRICULUM **DATA MINING**

Coordination: BEJAR TORRES, RAMON

Academic year 2023-24

Subject's general information

|   |  |               |                  |
|---|--|---------------|------------------|
| <b>Subject name</b>   | DATA MINING  |               |                  |
| <b>Code</b>   | 103089   |               |                  |
| <b>Semester</b>   | 1st Q(SEMESTER) CONTINUED EVALUATION   |               |                  |
| <b>Typology</b>   | <b>Degree</b>  | <b>Course</b> | <b>Character</b> |
|   | Master's Degree in Informatics Engineering                                       | 2             | OPTIONAL         |
|   |  |               | <b>Modality</b>  |
|   |  |               | Attendance-based |
| <b>Course number of credits (ECTS)</b>  | 6  |               |                  |
| <b>Type of activity, credits, and groups</b>                                    | <b>Activity type</b>   | PRALAB        | TEORIA           |
|   | <b>Number of credits</b>   | 3             | 3                |
|   | <b>Number of groups</b>  | 1             | 1                |
| <b>Coordination</b>   | BEJAR TORRES, RAMON  |               |                  |
| <b>Department</b>   | COMPUTER ENGINEERING AND DIGITAL DESIGN  |               |                  |
| <b>Teaching load distribution between lectures and independent student work</b> | 30% of the time are lectures (3 hours/week) and 70% is based on autonomous work. |               |                  |
| <b>Important information on data processing</b>                                 | Consult <a href="#">this link</a> for more information.                          |               |                  |
| <b>Language</b>   | English  |               |                  |

| Teaching staff      | E-mail addresses     | Credits taught by teacher | Office and hour of attention |
|---------------------|----------------------|---------------------------|------------------------------|
| BEJAR TORRES, RAMON | ramon.bejar@udl.cat  | 3                         |                              |
| MATEU PIÑOL, CARLOS | carles.mateu@udl.cat | 3                         |                              |

## Subject's extra information

To follow this subject, the student should have solid knowledge of structured, object and functional programming in python

## Learning objectives

1. To know the current tools for Data Cleaning and Data Analysis
2. To know the basics for the development of data-centric procedures using interactive programming tools
3. To know how to transform raw data from any source to consistent data for its analysis
4. To know how to implement and debug procedures for the transformation of massive data sets using Big Data approaches
5. To acquire a sceptic spirit in front of sets of data to incentive the exploratory analysis using computer science tools
6. To acquire the tools and knowledge for the descriptive analysis of potentially massive and intractable sets of data
7. To know and use the most basic data mining algorithms for discovering relevant features from big data sets
8. To know and use basic and advanced algorithm for machine learning suitable for big data applications
9. To know the basics of recommender systems

## Competences

### University of Lleida strategic competences

- **UdL2:** Command of a foreign language.

### Cross-disciplinary Competences EPS

- **EPS1:** Capacity of planning and organizing the personal work.
- **EPS3:** Capacity to convey information, ideas, problems and solutions to both a specialized and no specialized public.
- **EPS4:** Capacity to conceive, design and implement projects and/or contribute to new solutions, using engineering tools.

- **EPS5:** To be motivated for the quality and steady improvement.

## General Competences

- **CG4:** Capacity to mathematically model, calculate and simulate in technological companies and engineering centres, particularly with regard to research, development and innovation tasks in all fields related to computer engineering.
- **CG8:** Capacity to apply the knowledge acquired for solving problems in new and unfamiliar situations within broader and more multidisciplinary contexts, and to be capable of integrating this knowledge.

## Basic Competences

- **CB3:** Students are able to integrate knowledge and handle complexity, and formulate judgments based on information that was incomplete or limited, include reflecting on social and ethical.
- **CB4:** Students can communicate their conclusions -and the knowledge and rationale underpinning these, to specialist and non-specialist audiences clearly and unambiguously.

## Degree-specific competences

- **CE1:** Capacity to understand and apply advanced knowledge in high-performance computing and numerical or computational methods to problems of engineering.
- **CE4:** Capacity to model, design, define the architecture, implant, manage, operate, administer and keep applications, networks, systems, services and computer contents.

## Subject contents

1. Data mining and learning
  1. Mining frequent items/item sets and distinct elements
  2. Dimensionality reduction
  3. Linear and Logistic regression with SGD
  4. Naive Bayes classifiers
  5. Clustering
  6. Recommender systems
2. Neural networks and Deep Learning
  1. Introduction to neural networks
  2. Deep Neural Networks
  3. Convolutional Networks
  4. Recurrent Neural Networks
3. Control: Reinforcement Learning

## Methodology

Every week, each student will receive:

- Three hours of class attendance. Lectures are conducted by theoretical explanations accompanied by illustrative examples in the first part, finalizing with practical exercises in the second part. As a support material of the class, we will follow the slides or python notebooks of the course.
- Other support materials to follow the subject in a non-attendance way.

The evaluation is continuous throughout the semester and consists of four different parts:

- Two practices: Extending the bigdata application started at the previous subject "Massive data processing" with neural networks and big data tools.

- Two reports and oral presentations about integrating neural networks and data mining tools in their bigdata application.

**Due to the COVID-19 situation some classes will be non-presential, i.e. using videos and/or videoconferences.**

## Development plan

Weekly:

1. Mining frequent items/item sets.
2. Mining distinct elements.
3. Dimensionality reduction.
4. Linear and Logistic regression.
5. Naive Bayes classifiers.
6. Clustering - crisp and probabilistic.
7. Recommender systems.
8. Work on Data mining projects.
9. Work on Big data application.
10. Neural Networks.
11. Feed forward neural networks.
12. Deep Learning.
13. Convolutional Neural Networks.
14. Convolutional Neural Networks II.
15. Recursive Neural Networks.
16. Control: Reinforcement Learning.
17. Work on Deep Learning projects.
18. Work on Big Data Applications.
19. Final presentations of big data projects.

**Due to the COVID-19 situation some classes will be on-line, i.e., using videos and/or videoconferences.**

## Evaluation

The assessment for this course is based on continuous evaluation.

| <b>Evaluation activities</b>                    | <b>%</b> | <b>Dates</b> | <b>O/V (1)</b> | <b>I/G (2)</b> |
|---|----------|--------------|----------------|----------------|
| <i>Data mining application (3)</i>              | 40%      | End Sem.     | M              | G              |
| <i>Data mining application presentation (4)</i> | 10%      | End Sem.     | M              | G              |
| Deep learning exercises                         | 25%      | Middle Sem.  | M              | G              |
| Data mining exercises                           | 25%      | Middle Sem.  | M              | G              |

(1) Mandatory / Voluntary

(2) Individual / Group

(3) There will be a revision of the work performed by each member of the group

(4) Each member of the group has to participate equally in the presentation and answer questions from the evaluator

## Bibliography

- Wes McKinney. Python for Data Analysis: Data Wrangling with Pandas, NumPy, and Ipython. O'Really, 2012
- Holden Karau, Andy Konwinski, Patrick Wendell & Matei Zaharia. Learning Spark. O'Really, 2015
- Jure Leskovec, Anand Rajaraman & Jeffrey David Ullman. Mining of Massive Datasets. Cambridge University Press , 2014 (Find a copy at <http://www.mmds.org/>) It is also available at our Library, but only one user can borrow the book at the same time).