# DATA MINING – Project Guidelines

Alessio Cascione
alessio.cascione@phd.unipi.it

a.a. 2025/2026

# General Information

**Who?**

- Groups of min 2, max 3 people (ideally, heterogeneously distributed… e.g., 2 DS + 1 InfoUma);
- Insert your group at the following [link](link).

**What?**

- A Project:
  - Data Understanding & Preparation;
  - Clustering;
  - Classification and Regression;
  - Pattern Mining.
- A Report:
  - max 20 pages of text, including tables and figures;

only the report will be <u>evaluated</u>!

# General Information

**When?**

There are two options:

    **[1 week before the oral exam]:** complete project → NO BONUS POINTS

<div align="center"><b>OR</b></div>

    **[15/11/2025]:** mid-draft (Data Understanding & Preparation + Clustering) → 0.5 points

    **[31/12/2025]:** complete project → 0.5 points

So delivering the mid-draft and the complete project within the deadlines → 1 point

Deadlines can change, so you should refer to the main page of the course;

**Where?**

- send the report to **BOTH** alessio.cascione@phd.unipi.it & riccardo.guidotti@unipi.it

  - Object: [DM1 Project 25/26]

  - Report title: Project_Surname1_Surname2….pdf

**How?**

- Code language: Python (suggested);

- Report: written in LaTex (Overleaf) (suggested), or Office/OpenOffice…;

# The Project

**Board Games Dataset**: The dataset contains information regarding more than 20k board games rated by an online board game community. The information provided range from year of publication, game difficulty, number of players, information about games rank in specific categories and game characteristics/themes.

Links to the dataset will be provided on the main page of the course!

# The Report

## Structure

- Title page and index are not counted for the 20 pages limit;
- Only PDF are allowed, no doc, jupyter notebooks, python code;
- It is better to use font size higher than 9pt;
- Multiple columns are allowed;

## Content

- You must **justify every choice** (from the variables management to the parameters you tune);
- **Discuss every result**; even if some of them don't convince you, be fair and try to discuss the possible limitations (they can be imputed to the dataset, to an algorithm that does not fit with the dataset, etc…);
- **Plots** and **tables** without any comment are **useless**;
- Nice and **readable plots** make your analysis more understandable ;
- Even if you find a top configuration for your algorithm (e.g., k-means, k=5) you must **list which are the different parameters you tested and justify your choice;**

# The Tasks

- Data Understanding & Preparation;
- Clustering;
- Classification and Regression;
- Pattern Mining;

The next slides provide several analytical suggestions, but:

- You are allowed to organize the content of the complete project as you prefer;
- **There is a mandatory classification task, but you are allowed to identify further classification tasks as you prefer. The regression task can be chosen freely**;
- You are allowed to explore tools and methodologies not introduced during the lectures (e.g., feature selection methods, new plots, algorithms), but it is suggested to write me an email before;

# Data Understanding & Preparation (30 pts)

- Data Semantics
  - Introduce the variables with their meaning and characteristics;

- Distribution of the variables and statistics
  - Explore (single, pairs of…) variables quantitatively (e.g., statistics, distributions);

- Assessing data quality and possible variable transformation
  - Are there errors, outliers, missing values, semantic inconsistencies, etc?
  - Is it better to use for further modules transformed variables (e.g., log-transformated)?

- Pairwise correlations and eventual elimination of variables
  - Matrix correlation (analyse high correlated variables);

# Clustering (30 pts)

- Analysis by centroid-based methods
  - K-Means (mandatory), Bisecting K-Means (optional), X-Means (optional);
  - Choose the attributes, identify the best value of k, discuss the clusters.

- Analysis by density-based clustering
  - DBSCAN (mandatory), OPTICS (optional);
  - Choose the attributes, identify the best parameter configuration, discuss clusters.

- Analysis by hierarchical clustering
  - Choose the attributes, the distance function, analyze several dendrograms.

- Final discussion
  - Which is the *best* algorithm? Remember that *best* is studied w.r.t. several aggregate statistics, cluster distributions and w.r.t. the typology of algorithm used for that particular dataset;

# Classification and Regression (30 pts)

- Classification of **at least the Rating variable (mandatory)**:
  - by Decision Trees, KNN, Naive Bayes.

You should discuss the choice of the attributes and identify the best parameter configurations (e.g. gain criterion for trees, best k for KNN etc.).
Any other potential target variables for classification beside Rating can be chosen for extra analysis

- Regression - **single** and **multiple regression**:
  - Choose **one target** and **one independent variable** and solve a linear regression task
  - Choose **one target (the same for the previous task)** but this time consider **2+ independent variables** and solve with linear and at least 2 non linear approches

- Discussion
  - Evaluate the quantitative performance of the classification algorithms w.r.t. confusion matrix, accuracy, precision, recall, F1, ROC curve
  - Evaluate the quantitative performance of the regression algorithms w.r.t. MSE, $R^2$
  - Discuss some insight (e.g. try to interpret the tree(s))
  - Which is the *best* algorithm? *Best* can be studied w.r.t. the performance evaluation or other preferred point of view;

# Pattern Mining (30 pts)

- Frequent Pattern extraction
  - Using different values of support, etc;

- Discuss Frequent Pattern

  - Including qualitative and quantitative analysis, e.g., how the number of patterns w.r.t k min_sup changes;

- Association Rules extraction

  - Using different values of confidence, etc;

- Discuss Association rules

  - Including qualitative and quantitative analysis, e.g., how the number of rules w.r.t k min_conf changes, histograms of rules' confidence and lift;

- Exploit the most useful extracted rules

  - E.g., use them to replace missing values or to predict the target variable;

# Bonus & Other

- You can get 3 additional extra points in the final mark w.r.t. the following criteria:
    - Innovation (0.5 pts)
    - Experimentation (0.5 pts)
    - Performance (0.5 pts)
    - Appearance, Summary, Organization (0.5 pts)
    - Mid-draft and complete project within time (0,5 + 0,5 = 1 point)

- **Project Mark**: average of the previous modules + 3 bonus points;