# Individual Assignment 2: Data Complexity and Meta-Learning

**Submission Deadline: Nov 3**

## Objective

The purpose of this assignment is to explore the usefulness of meta-information when addressing machine learning tasks. **The deliverables of this assignment are a 2-page report with your insights and a jupyter notebook (or Python project) with your experiments. The assignment might be subjected to defense.**

## Instructions

In this individual assignment, you should identify an open-ended question to investigate using data complexity and meta-learning techniques. Your problem could focus on understanding classification difficulties, comparing algorithmic performance, exploring recently proposed measures, or designing new approaches. Some ideas are given below.

## Open-ended Projects

**A.** Explore a chosen dataset using *problexity*, *pymfe*, *pyhard* or others. Identify the main complexity factors associated with the data and reflect on their implications to *i)* classification algorithms with different learning paradigms or *ii)* preprocessing algorihtms. For *i)* you might select two or three classifiers with different paradigms (e.g., Decision Trees, KNN, SVM), while for *ii)* you can focus on resampling algorithms (e.g., cleaning approaches, oversampling). Verify whether your evaluation of meta-features matches the expected classification results. You can look into the research of Garcia et al. (2018).

**B.** Select a set of datasets for classification (e.g., 10, 15 datasets). Is it possible to link the meta-characteristics of these datasets to the obtained classification performance? For instance, you can select a single classifier and determine which hyperparameters work best for different complexity characteristics. You can also try to cluster datasets based on their meta-features and analyse whether they map onto clusters of good/bad behaviour. You can look into the research of Garcia et al. (2018).

**C.** Search for a recently proposed complexity measure/meta-feature that is yet to be implemented in the studied packages. Implement it and experiment with some datasets. What is the rationale behind the proposal? What (new or different) information, if any, does the new measure add to the state of the art?

**D.** One approach to determine the complexity of a dataset for classification refers to analyzing its data typology (Napierala et al. (2016)). Implement this logic and test it across several datasets. Can you find a way to order the datasets by their S/B/R/O percentages in a way that it matches the performance results?

**E.** Select a paper of your choosing and try to apply it to a different context or dataset. For instance, how would FAWOS (Salazar et al. (2021)) behave with other biased datasets? You can explore the technique across untested datasets to see whether the algorithm is feasible or if the results hold. Some well-known datasets can be found in Le Quy et al. (2022).

**F.** Can instance hardness measures hint at potential unfair situations? Lorena et al. (2024) conceptualize some ways in that instance hardness measures can help identify bias in sensitive groups. Produce a proof-of-concept with some experiments of your own. Is this a reasonable claim?

**Deliverables**

For this assignment, the following deliverables are required:

- A report (up to 2 pages -- using Template 1 or Template 2) summarizing your main findings. The report should include the following sections:
  - **Introduction:** Brief description of the problem definition and objectives;
  - **Methodology:** Explain the rationale of your approach and the steps taken;
  - **Results and Discussion:** In-depth discussion of your findings;
  - **Conclusions:** Summarize key findings, limitations, and learnings;
  - **References:** Cite all resources used in your work (e.g., books, research papers, websites, tools, packages). Use a consistent citation format.
- A Jupyter Notebook (or .py project) containing the complete code used for analysis, inline comments explaining important aspects and visualizations/data included in the report.

# Evaluation Criteria

Your submission will be evaluated based on the following criteria:

- **Creativity and Problem Definition**: Ability to define a clear and well-motivated question for investigation; originality in framing the approach to the analysis; creativity in the used tools.
- **Quality of Experiments**: Quality and thoroughness of the experiments designed to address the problem; appropriate use of meta-information to draw conclusions; correct implementation of methods, and clear motivation and explanation of approaches.
- **Reflections and Insights**: Depth of analysis and critical reflection on the obtained results; ability to interpret the implications of meta-information on model performance, preprocessing, or other aspects; discusion of limitations and challenges.
- **Clarity and Organization of Deliverables**: Writing and logical flow and structure of the report and the Jupyter Notebook (or project).
- **Bibliography**: Quality of the sources and tools and their relevance to the analysis.

# Bibliography

- Garcia, L. P., Lorena, A. C., de Souto, M. C., & Ho, T. K. (2018). Classifier recommendation using data complexity measures. In *2018 24th International Conference on Pattern Recognition (ICPR)* (pp. 874-879). IEEE.
- Le Quy, T., Roy, A., Iosifidis, V., Zhang, W., & Ntoutsi, E. (2022). A survey on datasets for fairness-aware machine learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, *12*(3), e1452.
- Napierala, K., & Stefanowski, J. (2016). Types of minority class examples and their influence on learning classifiers from imbalanced data. *Journal of Intelligent Information Systems*, *46*, 563-597.
- Salazar, T., Santos, M. S., Araújo, H., & Abreu, P. H. (2021). FAWOS: Fairness-aware oversampling algorithm based on distributions of sensitive attributes. *IEEE Access*, *9*, 81370-81379.
- Lorena, A. C., Paiva, P. Y., & Prudêncio, R. B. (2024). Trusting my predictions: on the value of Instance-Level analysis. ACM Computing Surveys, 56(7), 1-28.