

Inteligencia Artificial Explicable

Hands-on Project Unit 2

Contents

| | |
|----------------------------------|----------|
| Formal requirements | 1 |
| Submission | 1 |
| Deliverables | 1 |
| Phases | 1 |
| What to do? | 2 |
| Mandatory | 2 |
| Optional | 2 |
| How to do it? | 3 |
| Data | 3 |
| Code | 3 |
| Report structure | 3 |
| Presentation structure | 3 |
| Grading | 4 |

Formal requirements

Submission

- Strictly via Moodle.
- Upload before the deadline. Late submissions will not be accepted.
- **A single submission per group.**
- **Do not use Moodle comments** to deliver any sort of important information (e.g., group members). Include all relevant information in the deliverables.

Deliverables

1. A PDF file with a project report (maximum 15 pages). Must contain the names of group members.
2. All code used to carry out the analysis:
 1. Python: A notebook that must run, out of the box, on Google Colab.
 2. R: Script or RMarkdown file that must run, out of the box, in a typical R installation (e.g., please include code for installing required packages).
3. For the final delivery, the test set completed with predictions.

Phases

The project will be delivered in **two phases**. Regarding the first phase:

- use only the decision tree;
- corresponds to 10% of the project's grade (no minimum grade, and it is therefore not strictly necessary);

- deadline: 28/09/2025 at midnight.

The second phase corresponds to the delivery of the final and complete report and corresponds to 90% of the project's grade. For this second phase, you will be provided a test set, without labels, and you will fill-in those labels and upload the updated test set predictions to Moodle. Part of the evaluation grade depends on the accuracy on the test set. Failing to provide the test set in the adequate format will be considered as having baseline accuracy.

What to do?

Mandatory

You will train and interpret glass-box models. In particular, you will train at least two of the following classifiers:

- Logical (trees/lists/sets);
- Linear;
- Generalized additive models (Explainable boosting machine).

You may use any learning algorithm for the selected models, including standard greedy decision tree methods such as CART. You will also need to evaluate each of the classifiers on validation data. You will need to carry out both global and local interpretation; that is, interpret the learned models as well some of its predictions. In particular, regarding local explanations I recommend selecting 2-4 instances (at least one of each class) and explain their predictions with different models.

You will chose one of the models –the one that you consider the most adequate– and apply it to the provided test set, filling in the missing labels. You will submit the completed test set in the final submission.

Optional

The following are examples of additional tasks that can increase your grade:

- Try to increase the performance of the classifier on validation data, while using *training data only* to fit them – using validation data for training will lead to heavy penalties. This includes parameter tuning, preprocessing, etc.
- Try non-standard learning algorithms, such as those for learning optimized trees, or methods that we have not covered in class, such as RuleFit. Before using a method that we have not covered, check with the professor that the model is adequate.
- Evaluate black-box models such as random forest or gradient boosting, on the validation set, and compare to the performance of glass-box models.
- Discuss differences in performance among the different models.
- Critical discussion of the process and/or obtained results.
- ...

The above list is not exhaustive and you may try to improve the grade in a different way. If in doubt, check with the professor.

Note: some of these tasks may be time-consuming and yet produce no tangible results. For example, you might not be able to improve over the baseline method with parameter tuning, or the selected algorithm/package may be too slow. Before investing significant into into such a tasks try to improve the project in 'more straightforward' ways.

How to do it?

Data

The training set is contained in the **train.csv** file. It has eight columns, corresponding to eight variables. The first variable, 'RiskPerformance', is the target variable, to be predicted from the rest.

Code

You are free to use any Python or R package. You may combine Python and R. Some suggestions:

- scikit-learn;
- interpretml (<https://interpret.ml/docs/dt.html>);
- imodels (<https://github.com/csinva/imodels>);

Report structure

The report may be written in either Spanish or English. The structure of the report is as follows:

1. Cover page. Include a cover page with title, authors' names, email, course, and date. **Make sure to list all authors (group members).**
2. Introduction. Briefly explain the problem, the data, your overall approach and give overview of your results.
3. Method. Indicate your strategy and reasoning for fitting of each of the classifiers, justifying the choice of the algorithm, parameters and configurations. You may discuss why certain options might make sense and while others might not.
4. Results. Provide validation set accuracy results for each classifier.
5. Models. Interpret the learned models. This is the most important section.
6. Discussion: Discuss the models' performance and interpretability (subjective). Could some aspects be improved? For example, if you converted a tree to a rules list, could the list have been further simplified?
7. Conclusion. Provide conclusions regarding the project, your approach and your results. Discuss possible additional steps towards improvement.
8. References.

Results summary

Important: The Results section must begin with a table summarizing considered methods and their best prediction results. There will need be a row for each model, and a column for the best non-default parameter settings for that class as well as for the best training set and 5-fold CV accuracy obtained with those settings. For example:

| Classifier | Best parameters | Train accuracy | 5-fold CV accuracy |
|----------------------------|------------------|----------------|--------------------|
| ExplainableBoostingMachine | | 0.85 | 0.75 |
| RandomForestClassifier | n_estimators=100 | 0.85 | 0.75 |
| ... | ... | ... | ... |

Presentation structure

The presentations will be in class. Each presentation will last for 4 minutes, with 3 minutes for questions. I recommend having four slides: one summarizing what you have done and the results you have obtained, and then one for the interpretation of each of three different types of models (e.g. tree, linear, additive). Ideally, include some examples of local interpretation within the slides. If you do not have time to present everything you have done, choose for the presentation what you think is the most striking or has turned out the best. It is mandatory that both members of the group present.

The presentations will be made from my computer and that is why they have to be delivered the day before; *I will not accept presentations delivered the same morning*. The presentation has to be a *pdf file*, to avoid possible formatting problems.

Grading

Achieving a grade of 5 requires:

- Delivering all mandatory content, as specified above;
- Not-below-chance accuracy with each model;
- Absence of major conceptual or formal mistakes.

Secondary criteria include:

- The quality and clarity of the explanations, justifications, the writing and the presentation. Balancing clarity with rigor.
- Use of plots and other aids that help understanding the report.
- Extent of originality. Use of own data processing, methods, analysis.

Important: There will be penalties for not following submission instructions, such as exceeding the page limit, not providing the code, etc. Additional recommendations:

- Not everything that you do needs to appear in the report; keep the report concise and to the point. I can check the notebook to see all necessary details.
- Avoid plain copy-paste of code and output in the report; present the results in a clear and concise way. In particular, **do not deliver Jupyter or RMarkdown notebooks rendered as pdf**.
- Watch out with figures: for example, ensure the axis' scales are discernible.
- I encourage discussion of the achieved results, the motivation for trying some options and not others, and the reasoning behind the choices made, etc.