# Dropout Insight: Educational Risk Dashboard with Counterfactual Explanations

# User Manual

# Index

# Index of images

# 1. Introduction

This user manual aims to guide users in using the tool, offering a clear and simple description of its main features.

This document provides the necessary instructions for interacting with the tool, from initial access to executing the various options available.

The overall purpose is to help users install, run, and use the application effectively, as well as to provide solutions to any questions or problems that may arise during use.

This manual has been prepared voluntarily, considering that the tool requires a guide to help people learn how to use it properly.

The document is designed to support both primary users and those with more experience, facilitating understanding of the system through practical explanations and illustrative examples. In addition, recommendations for use and best practice guidelines are included to optimize the experience and make the most of the tool's capabilities.

In this way, users will have a comprehensive and accessible reference to resolve frequently asked questions, understand the operating logic of each module, and use the tool effectively and correctly.

With a clear and easy-to-understand approach, this manual is intended to be an additional, useful, and accessible resource for users, offering them confidence and guidance that will enable them to use the tool successfully.

## 2. Installation

In order to use the application correctly, it is necessary to follow a series of installation steps that ensure its proper functioning.

Since the entire project is coded in the Python programming language along with various associated libraries, it is important to take certain specific considerations into account during installation.

The tool is completely free to use and, for its construction, the Dash library has been chosen, which allows applications to be deployed locally using Flask.

As this is a project without financial investment, it has not been possible to have a free server with the power and space necessary to generate and store the various predictive models used during development. For this reason, the tool can only be run in a local environment. This choice not only responds to infrastructure limitations, but also gives the user complete control over their data and generated models, reinforcing the security and privacy of the process.

This procedure has been designed with a practical and simple approach, so that users can complete the installation independently, avoiding common errors and ensuring that the tool is ready for use in the shortest possible time.

Recommendations and best practices are also included, such as the use of virtual environments, which allow the tool to be kept organized and minimize conflicts with other programs installed on the computer.

The purpose of this section is, therefore, to provide the user with a comprehensive and orderly guide so that the installation can be carried out without complications.

The following sections describe step by step the actions necessary to carry out the installation, from preparing the environment to running the application for the first time.

## 2.1.    Prerequesites

Before beginning the installation process, you must check that your system meets the following requirements:

- Operating system: Windows 10/11 or recent Linux distributions (e.g., Ubuntu 20.04 or higher).

- Python: For the system to work properly, you must install version 3.9.0 exactly, due to the compatibility of the different libraries used.

- Pip: Updated Python package manager.

- Web browser: Google Chrome, Mozilla Firefox, Microsoft Edge, or equivalent updated browser.

- Internet connection: To download necessary dependencies and libraries.

## 2.2.    Code download

Once the prerequisites have been met, especially the installation of the required version of Python, the code can be downloaded. To do so, access the GitHub link via the following link:

GitHub

Once inside, simply press the green "code" button to view the different download options available.

- If you have Git installed, you can do "git clone marttamunzz/Dropout-Insight-Educational-Risk-Dashboard-with-Counterfactual-Explanations:" This project develops an experimental dashboard to predict and explain the risk of academic dropout. It combines AutoML for model training with SHAP and individual and group counterfactuals to make predictions transparent and interpretable." In the terminal, download the tool to the directory where you are located.

- On the other hand, if you do not have Git installed, you can download it as a .zip file.
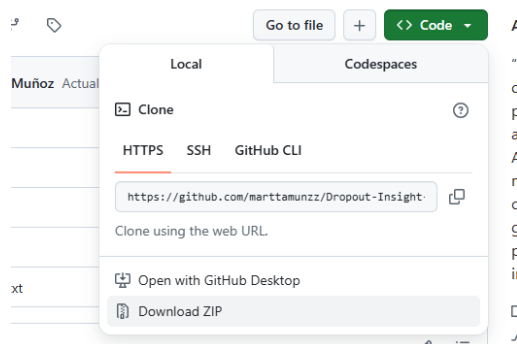


Image 1: zip download option

## 2.3.    Virtual enviroment creation

To ensure proper functioning and avoid conflicts with other libraries already installed on the system, it is recommended to perform the installation within a Python virtual environment, which allows project dependencies to be isolated so that they do not interfere with those of other programs or applications on the same computer.

It is important to note that the virtual environment only needs to be created the first time. Once created, each time you want to work with the tool, you only need to activate and deactivate the corresponding environment.

The steps required to create and activate a virtual environment are described below:

- Open the terminal or command line
- Go to the project's root folder
- Create a virtual environment using the following commands:

**python -m venv venv**

In this case, the environment will be created in a folder called venv, although you can assign another name if you wish

- Activate the newly created virtual environment:

Windows:

**venv\Scripts\activate**

Linux/MacOS:

**source venv/bin/activate**

Once activated, the name of the environment will appear in parentheses at the beginning of the command line, indicating that the virtual environment is in use:

**(venv) C:\User\Project>**

- When you want to end the work session, you must exit the virtual environment as follows:

**deactivate**

## 2.4.   Installing the required Python libraries

Once the virtual environment has been created and activated, the next step is to install the libraries that the tool requires to function properly. These libraries are all listed in the "requirements.txt" file, which contains the list of dependencies with their corresponding versions.

An internet connection is not necessary for this process, as the libraries must be downloaded from the official repository.

To complete the installation, follow these steps:

- Check that the virtual environment is activated: The name of the environment should appear at the beginning of the command line.

- Run the following command in the project's root folder:

**pip install -r requirements.txt**

This command will automatically download and install all the necessary libraries from the official Python repository (PyPI)

- Verify the installation of the libraries by running:

**pip list**

This will display a list of all the libraries installed in the virtual environment, along with their versions.

In this project, the most important libraries that we need to have installed, without which the tool cannot be used correctly, are the following:

**Dash (2.0.3)**

**Flask (2.3.2) + Werkzeug (2.3.4) + waitress (2.1.2),**

**pandas (2.3.1)**

**numpy (1.23.0)**

**scikit-learn (1.2.2)**

**mljar-supervised (0.11.5)**

**lightgbm (3.3.5), xgboost (1.7.5), catboost (1.2)**

**explainerdashboard (0.4.2.2)**

**shap (0.41.0)**

**dice-ml (0.12)**

## 2.5.    Uninstallation

If the user no longer wishes to use the tool, it can be easily uninstalled. This process frees up space on the computer and removes the dependencies associated with the project.

If a virtual environment has been created, simply delete the project folder along with the folder corresponding to the environment (usually called venv or env).

If the application has been installed without a virtual environment, in addition to deleting the folder containing the tool's source code, it is recommended to manually delete the libraries that were installed on the system. To do this, access the Python installation path and, within the Lib folder, delete those libraries that are no longer necessary, thus avoiding unnecessary space usage.

# 3. Interface

The application interface has been designed to be clear, intuitive, and easy to use, even for users without advanced technical knowledge.

This section describes the different visual elements and interactive components that make up the interface, as well as their functionalities. In addition, all graphics include a camera icon in the upper right corner, which allows them to be easily downloaded in PNG format for later use or analysis.

## 3.1.   Main screen

Since this is an application that only works locally, to open it we must access the src folder in the terminal and then type the following in the command line:

python index.py

Once this is done, a tab will automatically open in the browser with the address:

**<u>Dropout Insight: Educational Risk Dashboard with Counterfactual Explanations</u>**



Image 2: Index screen

When running the tool, the user accesses the home screen, which displays the initial instructions, the option to load the dataset, and the possibility of downloading a CSV file template that must be uploaded.

This interface is designed to guide the user in the correct preparation of the data before starting the analysis.

Before loading the dataset, a series of conditions that must be met to ensure that the tool functions logically are indicated and displayed on the interface so that the user is aware of them.

In the "load dataset" section, the user has an interactive box for loading the data, which allows two options: either drag and drop the file directly onto the box, or manually select the file from the system's file explorer.

Once the data file is uploaded in CSV format, a table is displayed on the screen showing all the instances in the file and a checklist to verify that all the specified requirements are met.



Image 3: Data uploaded in a table

If a dataset is uploaded that does not meet any of these requirements, the tool displays an error message indicating the problem and what needs to be done, as well as highlighting in red the item on the checklist that has not been met.
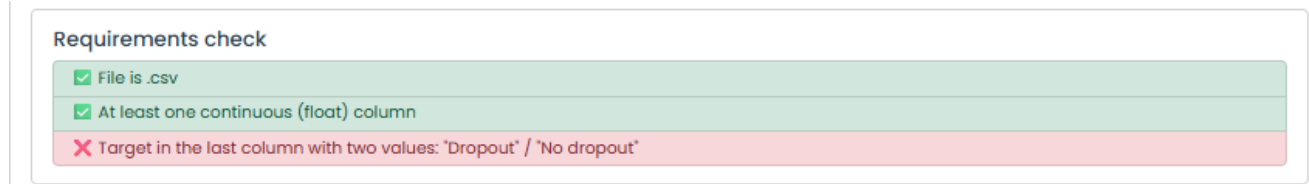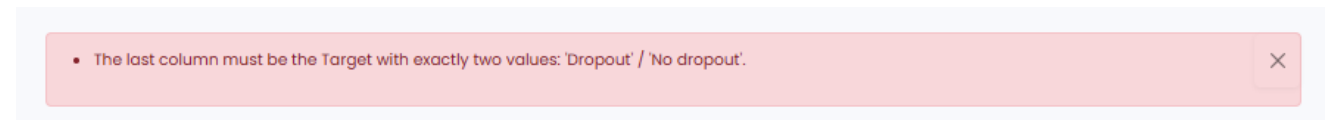
Image 4: checklist with fail



Image 5: Wrong file error message

Once the correct dataset has been loaded, the user must press the "start" button for the tool to proceed to display the dashboard hub.
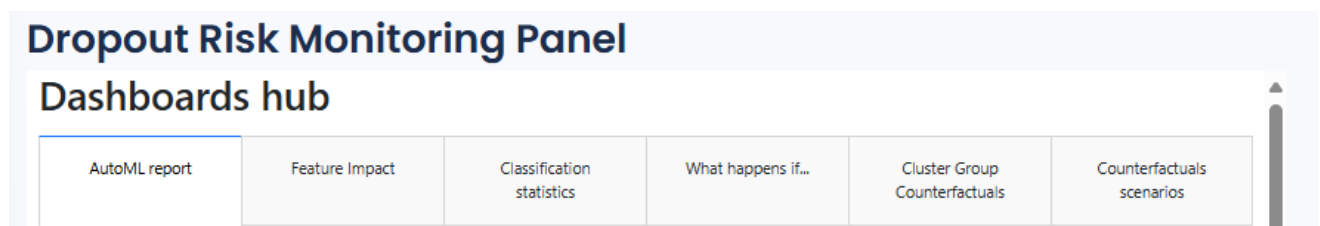


Image 6: Dashboard's hub

Once loaded, the user can choose which tab to navigate to, although it is recommended to follow the order in which they are presented:

## 3.2. Dashboard hub loading screen

When running the tool, the user accesses the Hub's main screen, which is the entry point for all the system's functionalities.

This interface has been designed to offer a clear, intuitive, and accessible experience for both technically knowledgeable users and those who are unfamiliar with machine learning tools.

From this screen, users can navigate the analysis modules included in the tool, such as viewing predictions, interpreting the importance of variables, or performing counterfactual analysis at both the individual and group levels.

The interface design of this tab responds to a user-centered approach in which the main elements are clearly highlighted, and the available actions are organized hierarchically, making the workflow easy to understand. In addition, the inclusion of interactive elements improves usability and makes the tool more personalized.

### 3.2.1. AutoML report tab

This tab is one of the most important sections of the tool, as it automatically collects the results obtained during the model training process, offering the user a complete overview of the performance of the different models generated.

This section includes different blocks of information:

- **Model leaderboard**: a ranking of the trained models is displayed according to the evaluation metric used (logloss). For each model, data such as the type of algorithm, the value obtained from the metric, and the training time used are presented. This allows you to quickly compare which configurations have delivered the best results.

**AutoML Leaderboard**

| Best model | name | model_type | metric_type | metric_value | train_time |
|------------|------|------------|-------------|--------------|------------|
| | 1_DecisionTree | Decision Tree | logloss | 0.360425 | 23.77 |
| | 2_DecisionTree | Decision Tree | logloss | 0.349445 | 11.56 |
| | 3_DecisionTree | Decision Tree | logloss | 0.349445 | 12.5 |
| | 4_Default_Xgboost | Xgboost | logloss | 0.312559 | 13 |
| | 5_Default_RandomForest | Random Forest | logloss | 0.335174 | 10.27 |
| | 6_Xgboost | Xgboost | logloss | 0.307738 | 12.09 |
| | 8_RandomForest | Random Forest | logloss | 0.336298 | 11.62 |
| | 7_Xgboost | Xgboost | logloss | 0.323101 | 12.25 |
| | 9_RandomForest | Random Forest | logloss | 0.321659 | 17.27 |
| the best | Ensemble | Ensemble | logloss | 0.305976 | 3.75 |

Image 7: AutoML leaderboard

- **Performance graphs:**

1. <u>AutoML performance</u>: a scatter plot representing the evolution of the metric value throughout the different training iterations.
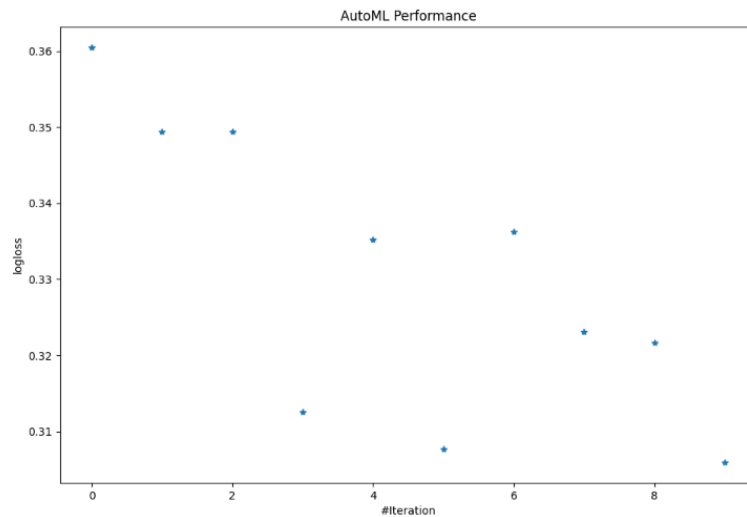


Image 8: Gráfico AutoML performance

2. <u>Performance Boxplot</u>: a box plot that summarizes the distribution of the values <u>AutoML</u> obtained, allowing for analysis of the variability and consistency of the models
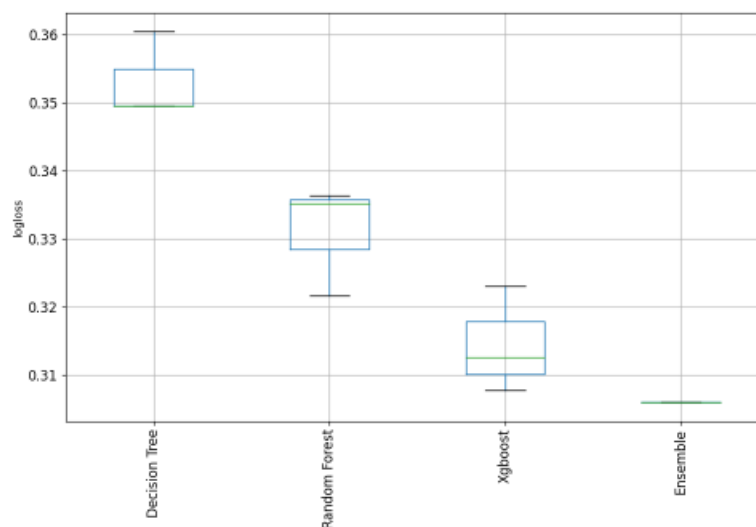


Image 9: AutoML Performance Boxplot

3. <u>Features importance</u>: A heat map showing the variables with the greatest influence on the predictions made by the best models. This block is particularly useful for interpreting which factors are most decisive in the risk of academic dropout.
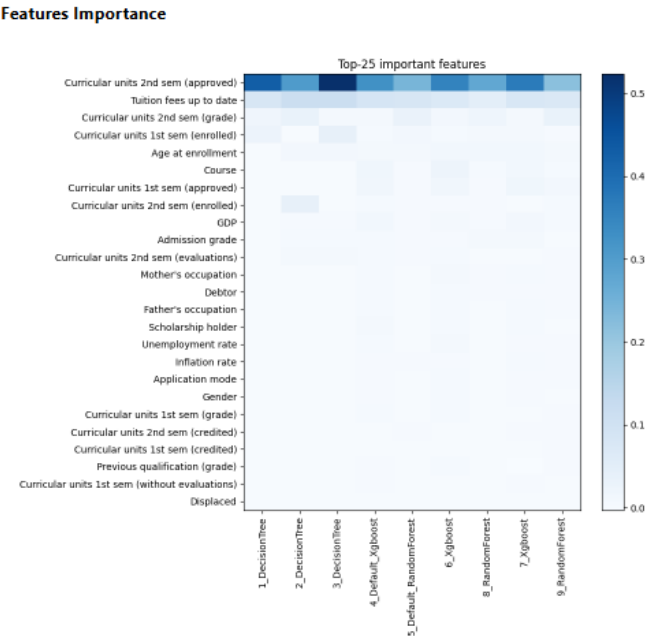


Image 10: Features importance

4. <u>Model correlation</u>: A graph that measures the similarity in the behavior of different trained models. With this information, you can evaluate the diversity between models and select more complementary combinations.
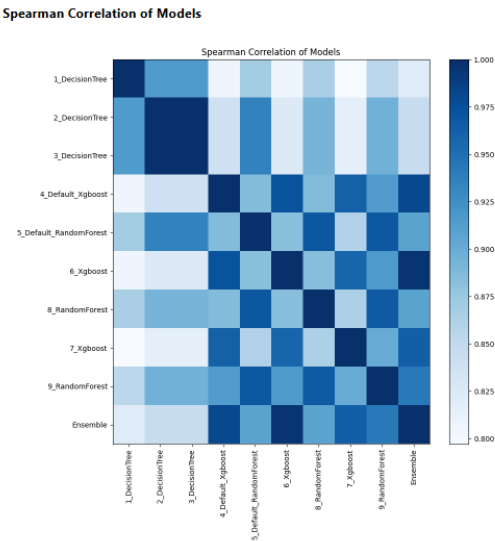


Image 11: Correlación entre modelos

### 3.2.2. Feature impact tab

The feature impact tab is intended for interpreting the models used to calculate SHAP (Shapley Additive exPLanations) values. This section allows users to understand which variables have the greatest influence on predicting academic dropout risk and how they impact the results.

This tab includes three main types of graphical visualizations:

- **Feature impact**: Shows using a bar chart, the average impact of each variable on the prediction. The calculation is performed using the average absolute value of the SHAPs, so that the user can directly identify which characteristics have the greatest weight in the model.
  In addition, there is a drop-down menu called depth, which allows you to select the number of variables you want to display in the graph.
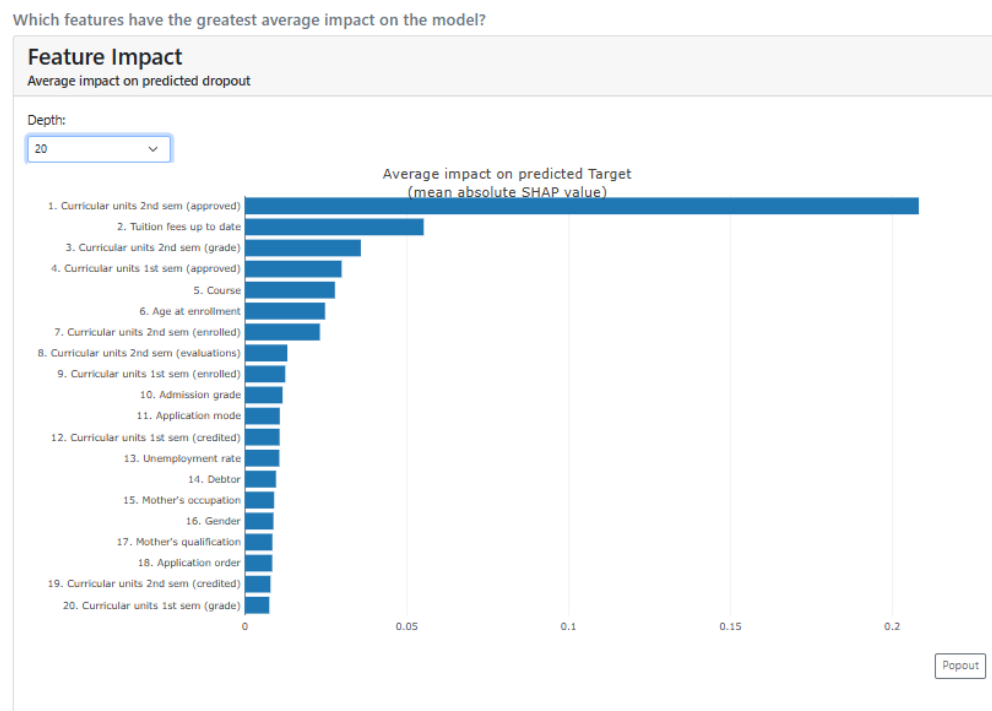


Image 12: Feature impact graph

- **SHAP Summary:** provides a more detailed view of the influence of each variable. This graph combines the effect (positive or negative) with the magnitude of the values of each characteristic in the prediction, using a color scheme that represents the intensity of the variable values (red = high impact value, blue = low impact value). Thanks to this visualization, global patterns can be detected and the direction in which each characteristic affects the prediction can be understood.In addition, in the "Summary type" drop-down menu, the graph includes the option to display in aggregate mode, which summarizes the information in a manner equivalent to the Feature Impact graph, allowing for a more concise interpretation.
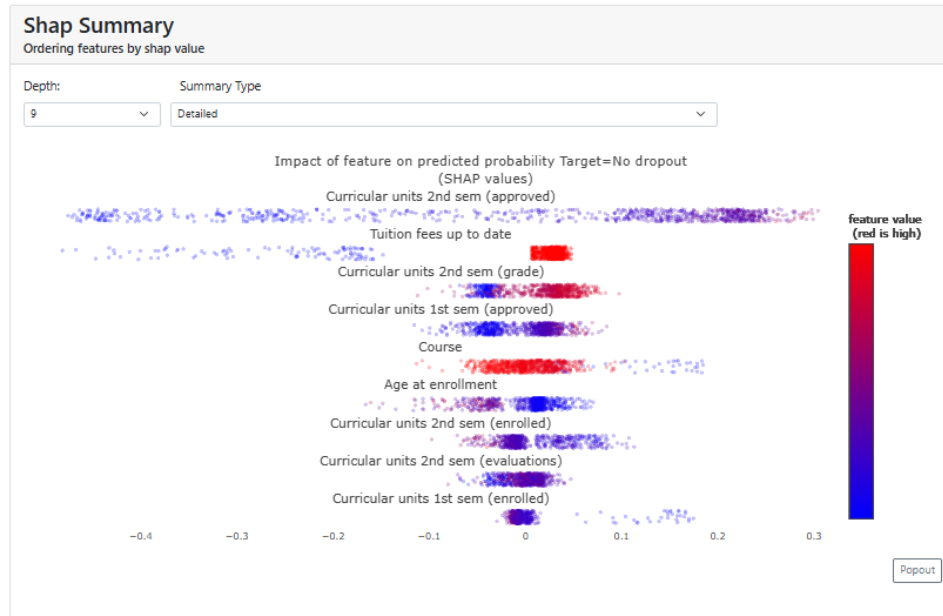
Image 13: SHAP summary per colour distribution

- **SHAP Dependence:** allows you to analyze in detail the relationship between the specific value of a variable and its impact on the prediction. The graph incorporates several drop-down menus that increase its flexibility:

    o Feature: selects the main variable on which the dependency is analyzed.
    o Color feature: allows you to color the points according to another variable, thus facilitating the detection of interactions between characteristics.
    o Fake_studentid: allows you to highlight a specific student within the graph to observe their individual behavior.

Thanks to these options, more complex patterns can be identified. For example, it is possible to observe how the number of subjects passed influences the probability of dropping out, conditioned by whether or not tuition fees have been paid.



Image 14: SHAP dependence relation between 2 features and target

### 3.2.3. Classification statistics tab

This tab provides users with a set of metrics and visualizations that allow them to evaluate the performance of the trained model. Its purpose is to provide a clear and quantitative overview of the system's ability to predict the risk of academic dropout.

The main elements included in this tab are:

- **Performance metrics**: these present classic evaluation indicators in classification problems such as accuracy (overall precision) and recall (sensitivity). F1-score (balance between precision and recall), ROC AUC and PR AUC (areas under the ROC and PRM curves that measure the model's discriminatory capacity), and log loss (logarithmic loss, which reflects the probabilistic quality of the predictions). These metrics allow for an overall assessment of the model's reliability.



Image 15: Statistic performance results

- **Confusion matrix:** represents the proportion of correct predictions, distinguishing between false positives, false negatives, true positives, and true negatives. Thanks to this

matrix, the user can identify not only the percentage of correct predictions, but also the most frequent types of errors.
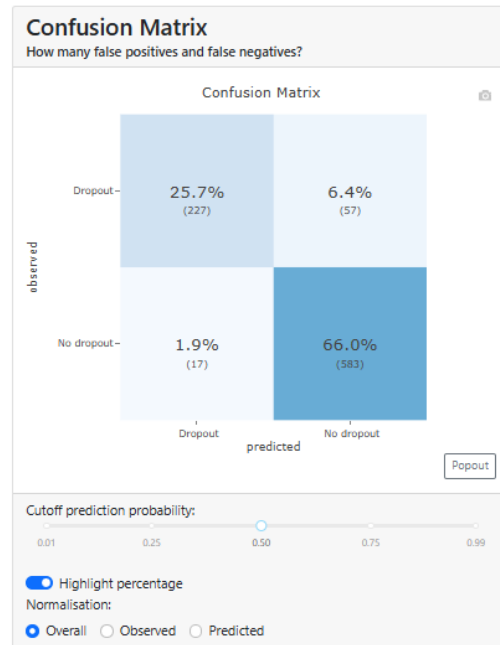


Image 16: Confusion matrix

- **Curva ROC (ROC AUC Plot):** Shows the relationship between the true positive rate and the false positive rate for different decision thresholds. The area under the AUC curve reflects the model's ability to discriminate between students at risk of dropping out and those who are not:
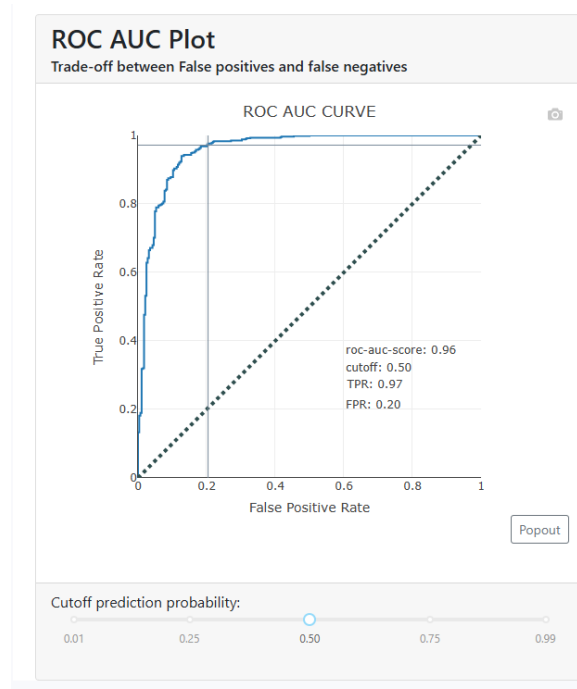


Image 17: ROC Curve

- **Classification plot:** Visualizes the distribution of predictions based on the selected cutoff threshold. This graph allows you to understand what proportion of cases are classified as dropouts versus non-dropouts and how these proportions vary if the probability threshold is modified.



Image 18: Classification Graph

### 3.2.4. What Happens if tab

This tab is designed for individualized analysis of predictions, allowing you to explore how the values of each characteristic influence the risk of academic dropout for a specific student. Its purpose is to answer the question, "What would happen if I changed one of the student's conditions?", offering an interactive simulation environment.

This tab is organized into several functional blocks:

- **Timer**: Allows you to start and stop a clock using the "start" and "stop" buttons. The objective is to measure the time spent analyzing or interacting with a specific case. This

functionality can be useful in practical or experimental sessions, where you want to record the time spent exploring different modifications in a student's variables.
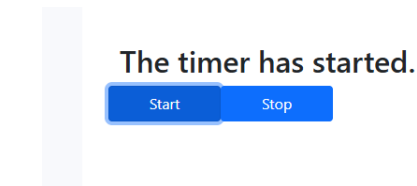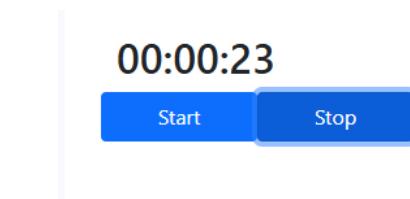


Image 19: Start of the clock



Image 20: End of the clock

- **Selecting a student:** The user can choose a specific student from the list using a drop-down menu or select one at random, allowing them to customize the analysis and focus on specific cases.



Image 21: Select a student option

- **Contributions plot:** Graph showing how each characteristic contributes to the model's final prediction. The green bars represent factors that reduce the probability of dropout, while the red bars represent those that increase it, allowing the user to clearly identify which variables have a decisive influence on the situation of a specific student.

Image 22: Contributions plot

- **Individual prediction:** It presents the estimated probability of dropout and non-dropout in numerical and graphical format (pie chart). This allows you to immediately see how each change in the variables would modify the final prediction.
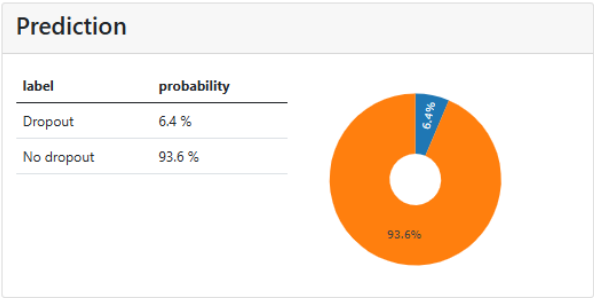


Ilustración 23: Individual predictions plot

- **Feature Input:** Interactive module that allows you to modify student variable values and observe how the probability of dropout changes in real time. This component is essential for simulating hypothetical scenarios, as it makes it easy to check how risk varies if the student modifies any of their variables.

Image 24: Table for modifying the features

The design of this tab responds to an approach focused on local interpretability by giving the user the possibility to manually experiment with the factors that influence the prediction. This makes it a complement to counterfactuals so that, while counterfactuals propose automatic solutions generated by the system that guarantee a high percentage change from one prediction to another, the What happens if... tab encourages interactive exploration and intuitive understanding of how each prediction is constructed.

### 3.2.5. Cluster Group Counterfactuals tab

The cluster group counterfactuals tab is geared toward collective analysis of academic dropout risk by generating counterfactuals at the group level. Its objective is to identify common patterns within sets of students with similar characteristics and propose modifications that could improve predictions across the board.

The system applies clustering techniques to form groups of students who share a similar academic and risk profile. Once the clusters have been identified, a representative student is selected from each group on which a counterfactual is generated. Subsequently, the modifications recommended in that case are extrapolated to the rest of the cluster members, showing how such interventions could improve the predictions for the group as a whole.

Before generating these counterfactuals, the tool offers a configuration panel where the user can define the cluster creation method and select which one to analyze.

There are two options:

- **Automatic**: the system automatically generates the optimal number of clusters based on internal criteria. This option is designed for users who prefer quick and assisted configuration.

- **Manual**: allows the user to manually set the number of clusters into which they want to divide the student set. This option offers greater flexibility and is useful when comparing the behavior of the model with different partitions. In this option, it is mandatory to indicate the number of clusters (k) to be generated.

Once the cluster creation option has been selected, click on the "create clusters" button and the system will generate them and display them in the form of a bar chart.
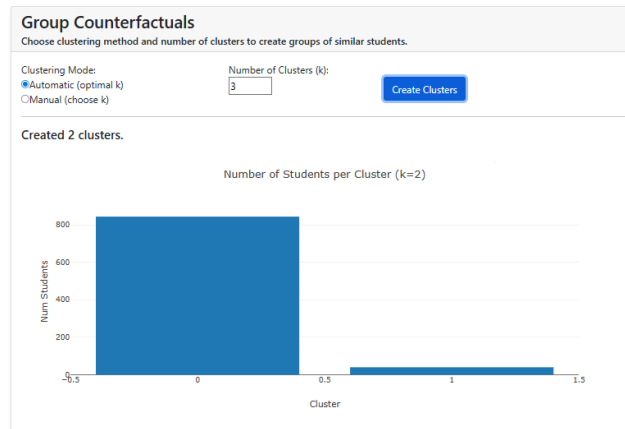


Image 25: Cluster creation

Once the clusters have been created, the tool displays a drop-down menu with the option to select the cluster group you want to analyze. Once you have selected one of the clusters, click on the "Generate Group Counterfactuals" button to display the different graphs for the proposed counterfactual.
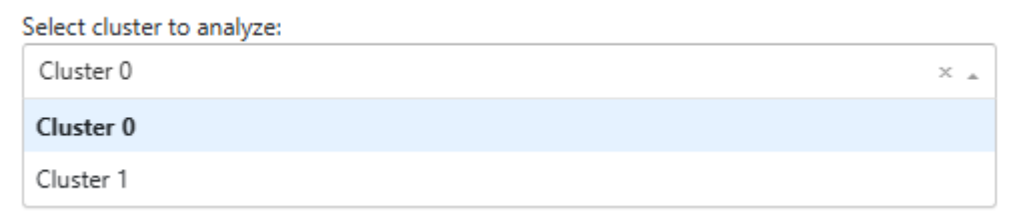


Image 26: Cluster selection

Once the group counterfactuals have been selected and generated, the tab offers a set of visualizations and metrics that allow you to evaluate their impact on the analyzed cluster.

- **Comparative tables**: first, two tables are displayed

  • Counterfactual table: shows the modified values for the representative student, highlighted in green, indicating the minimum changes needed to improve the prediction

  • Original student table: this shows the initial values for the same student, without modifications. This table serves as a reference to clearly visualize which attributes have been modified in the counterfactual.

Image 27: Counterfactual vs original data

- **Distribution of predictions:** A bar chart comparing the proportion of students classified as dropouts and non-dropouts before and after applying the counterfactual. The original scenario is shown in blue and the result after applying the changes is shown in orange.
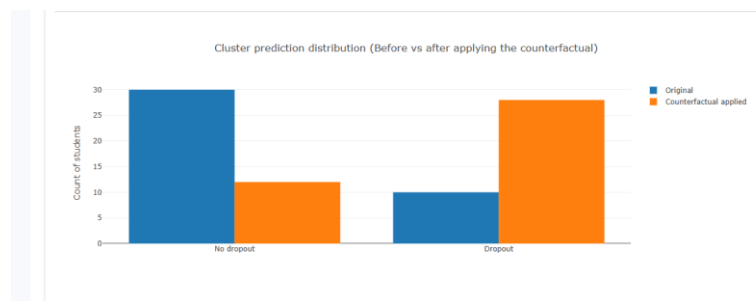


Image 28: Distribution of the prediction before and after the counterfactual

- **Number of modified features:** A histogram is included showing how many variables had to be modified, on average, to adapt each student to the group counterfactual, allowing the level of effort required to be assessed.
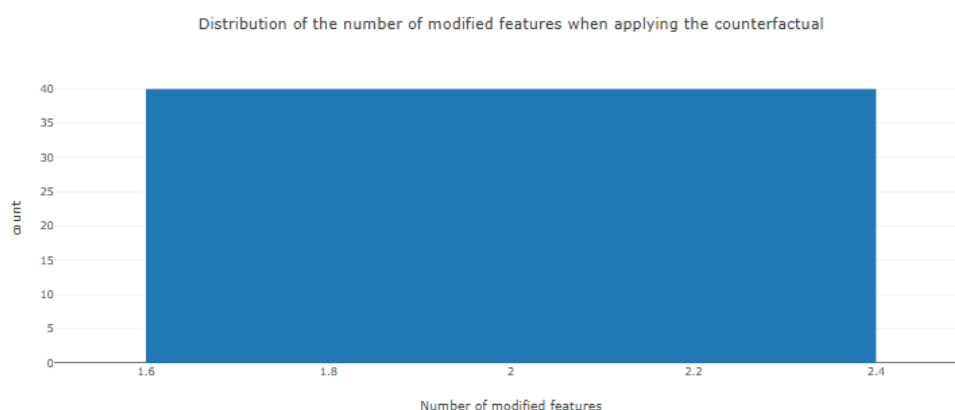


Image 29: Histogram of number of modifies features

- **Probability gain:** A box plot is included that represents the increase in the probability of success obtained after applying the counterfactual to the cluster. The greater this gain, the more effective the proposed scenario is.
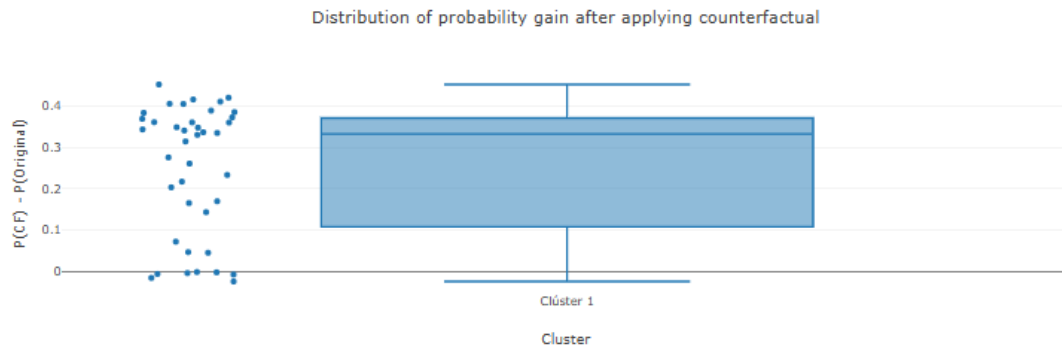
Image 30: Boxplot

- **Cluster summary:** Finally, a summary table summarizes the most relevant information, such as the number of students in the cluster, the original dropout rate compared to the counterfactual rate, the average number of modified variables, and the most modified variables within the group.



Image 31: Summary table of the group counterfactual

### 3.2.6. Counterfactuals scenarios tab

The individual counterfactual scenarios tab is designed to generate and visualize counterfactuals at the individual level, showing what minimal changes in the variables of a specific student would allow the dropout prediction to be changed to a scenario of academic retention.

The user selects a specific student, indicates the number of scenarios they want to be created, and the system automatically calculates a set of counterfactual scenarios, each of which represents an alternative version of the student, in which only some variables are modified in order to reduce the probability of dropout.



Image 32: Selection of student and numer of counterfactuals

Once the student has been indicated and the number of clusters selected, click on the generate scenarios button and the following components will automatically be displayed:

- Comparative tables:

    o Table of original values: with the values of the variables unchanged.
    o Table with the generated counterfactuals: where each scenario represents an alternative combination of values that would alter the initial prediction.



**Original**

| Marital status | Application mode | Application order | Course | Daytime/evening attendance | Previous qualification | Previous qualification (grade |
|---|---|---|---|---|---|---|
| 1 | 1 | 1 | 9130 | 1 | 1 | 1! |

**Counterfactuals**

| Marital status | Application mode | Application order | Course | Daytime/evening attendance | Previous qualification | Previous qualification (grade |
|---|---|---|---|---|---|---|
| 1 | 39 | 1 | 9130 | 1 | 1 | 1! |
| 1 | 1 | 1 | 9130 | 1 | 1 | 1! |
| 1 | 1 | 1 | 9130 | 1 | 1 | 1! |
| 1 | 1 | 1 | 9130 | 1 | 1 | 1! |
| 1 | 1 | 1 | 9130 | 1 | 1 | 1! |

Image 33: Original features vs counterfactuals

- **Probability visualization:** A bar chart is included that compares the probability of student dropout in the original scenario with the probabilities obtained in each counterfactual. This allows the user to identify which counterfactual scenarios offer the greatest reduction in risk.
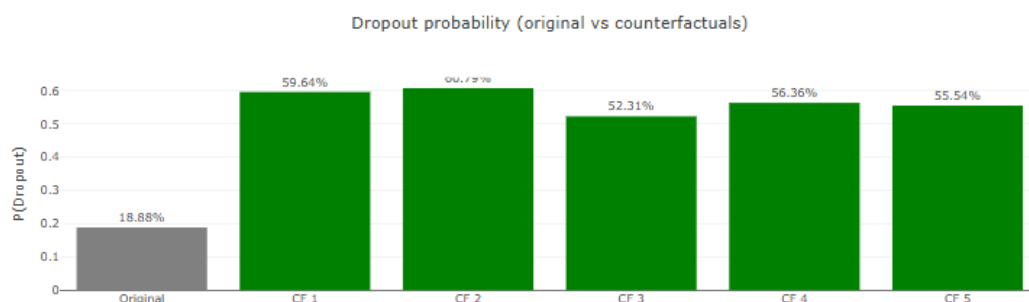


Image 34: Bar graph of dropout probability (original vs counterfactuals)

- **Radar Chart (original vs contrafactuals):** Graphically represents the differences between the student's original profile and the generated counterfactuals. Each axis corresponds to a modified variable, and the colored lines show how its value changes in each alternative scenario. This allows for an intuitive visualization of which variables have undergone modifications and to what extent.
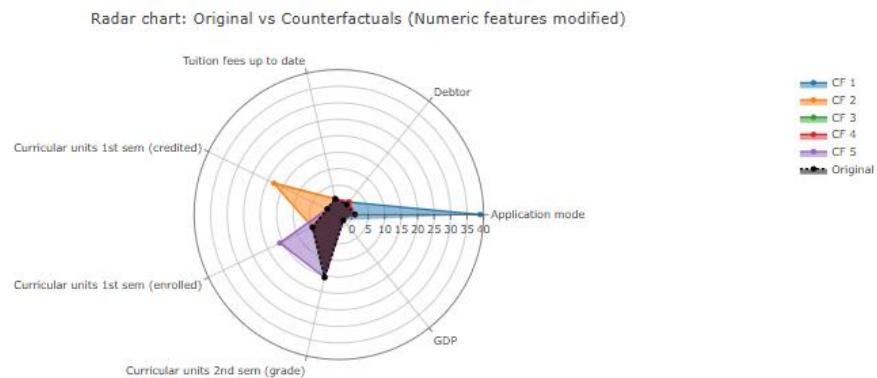
Image 35: Radar chart of the modified features

- **Changes in features table**: Summarizes the original values and the modified values in each counterfactual for the most relevant variables. The highlighted cells allow you to immediately identify which attributes have been modified in each scenario and with what exact value.

**Changes in features**

| Variable | Original | CF 1 | CF 2 | CF 3 | CF 4 | CF 5 |
|---|---|---|---|---|---|---|
| Application mode | 1 | 39 | 1 | 1 | 1 | 1 |
| Debtor | 0 | 1 | 1 | 0 | 1 | 0 |
| Tuition fees up to date | 1 | 1 | 1 | 0 | 1 | 0 |
| Curricular units 1st sem (credited) | 0 | 0 | 18 | 0 | 0 | 0 |
| Curricular units 1st sem (enrolled) | 5 | 5 | 5 | 5 | 5 | 16 |
| Curricular units 2nd sem (grade) | 15.6 | 15.6 | 15.6 | 15.6 | 12.4 | 15.6 |
| GDP | -1.7 | -1.7 | -1.7 | -2.86 | -1.7 | -1.7 |

Image 36: Table of the modifications suggested

- **Cumulative impact on features**: Finally, a bar chart shows the total contribution of each modified variable to the change in the probability of dropout. This section allows us to identify the most influential variables in the set of counterfactual scenarios, providing an overview of the impact of each characteristic.
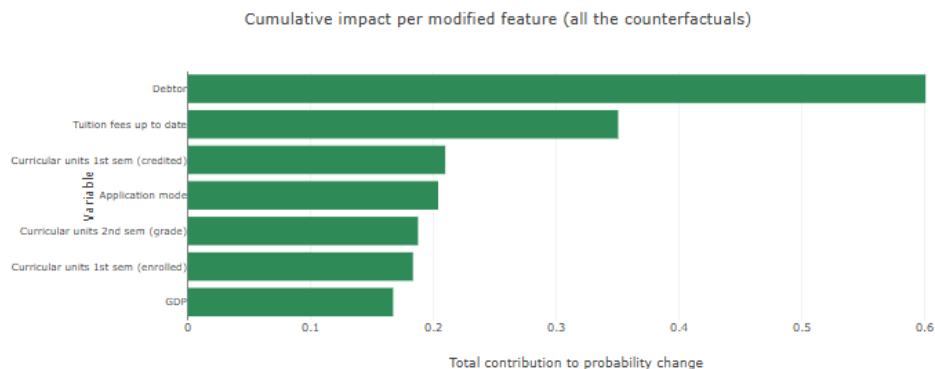


Image 37: Cumulative impact bar graph