# Data Science Pipeline Development

**Problem Introduction:**

In this challenge, your mission is to develop a robust pipeline capable of handling a classification task aimed at predicting a binary target variable (y-class). The pipeline will take datasets as input and output models that meet the criteria specified in the "Performance Objectives" section.

The focus on a pipeline arises from the requirement to develop multiple models for different datasets, which share identical columns and features but still require distinct models.

You are provided with two separate datasets to test the pipeline's robustness, ensuring it can independently handle each dataset. Combining the datasets is not permitted, as the goal is to generate different models for each dataset.

**Datasets Overview:**

**Dataset Composition:** The dataset contains approximately 50,000 rows, each consisting of 43 numeric feature columns (X1 to X43). The target variable, labeled 'True' or 'False,' is represented in the X44 column, with 'False' being the predominant class.

**Data Quality:** The data quality is pristine—there are no missing or incorrect values, as all invalid data were meticulously excluded during dataset creation.

**Feature Characteristics:** The features are all numeric, with some ranging from 0 to 100 (percentage scale) and others having no range limit. The relationships between these features and the target variable are likely non-linear.

**Data Order Importance:** The order of data is crucial as it reflects the sequence of real events. Therefore, it is essential to maintain this order in your analysis. The data used for testing should always follow the data used for training in sequence to preserve the integrity of the dataset's chronological order. Please do not shuffle the data.

**Performance Objectives:**

The ultimate goal of this project is to achieve, for each dataset, a precision of at least 60%, with the pipeline making a minimum of 50 true positive decisions for each 10% segment of the data in testing. At least 10% of the data should be allocated to the test set for each dataset, though allocating at least 20% is preferable for a more robust evaluation.

**Step-by-Step Goals:**

**Step 1: Initial Benchmarking**
The first objective is to demonstrate that the pipeline performs better than random decision-making. In the test set (which should be at least 10% of the data as mentioned), the initial target is to achieve at least 0.4 precision and 0.1 recall for each dataset. The target class is imbalanced, with the 'False' class making up roughly 75% of the data (though the exact ratio varies across datasets). Precision values around 0.25 indicate the model is performing at a random decision-making level. Surpassing this threshold shows that the model is making meaningful predictions, which is why 0.4 precision has been chosen as the first improvement target. This reflects clear progress beyond random performance.

**Step 2: Fine-Tuning for Intermediate Performance**
After the initial milestone is achieved, the goal is to refine the pipeline to improve performance further. The target is to achieve precision beyond 40% and recall of at least 0.1 (preferably higher, though this is not mandatory) for each dataset. The focus at this stage should be on enhancing the model's ability to make true positive decisions while maintaining a balance between precision and recall.

**Step 3: Achieving the Final Performance Goal**
The final objective is to reach a precision of at least 60% for each dataset. The model should consistently make at least 50 true positive decisions for each 10% segment of testing data for each dataset. The pipeline must remain robust and adaptable to different datasets, effectively handling the imbalanced class distribution without overfitting.

**Important Note**
A high recall (e.g., 0.95) combined with a precision value close to 0.25 (reflecting the class imbalance) suggests the model is making random decisions and is not reliable. Therefore, the pipeline must be carefully fine-tuned to improve both precision and recall, ensuring the model makes meaningful and non-random predictions.

**Suggestions for the Project:**

**1. Understand the Data and the Problem**
Before diving into coding, ensure you fully understand the problem. Pay attention to the imbalance in the target class and make sure you are clear on the performance objectives. Precision is the key metric in this project, but balancing it with recall is also important to meet the mission goals.

**2. Perform Exploratory Analysis Before Proceeding**
There is a high likelihood that the relationships between the features and the target variable are non-linear, and we are not currently aware of the feature importance. We strongly recommend that you conduct thorough exploratory analysis before proceeding with modeling. This exploration will help you better understand the data and guide decisions on feature selection, engineering, and the choice of models. It's very likely that feature selection and engineering will play a key role in enhancing model performance.

**3. Prioritize Precision Over Recall in Scoring**
Metrics such as the F1-score may not be ideal for this task, as precision is more important than recall in this case. Be mindful of this when selecting evaluation metrics, ensuring that the solution aligns with the project's precision-focused requirements.

**4. Handle Class Imbalance Correctly**
If class imbalance in the target variable needs to be addressed, ensure that balancing techniques are applied **only** on the training set. The test set must remain unaltered to reflect real-world conditions and provide accurate performance metrics.

**Process of Delivery and Payment:**

**1. Initial Understanding and Price Evaluation**
The first step is to carefully review this document, ensuring you fully understand the mission, including our data privacy, non-disclosure, and data security requirements (outlined in the next sections). Once you've reviewed and accepted these terms, the datasets will be provided to you. After that, we expect you to provide a price evaluation for the project. Once we agree on a price and timeline, the payment process will proceed as follows:

**2. Payment Structure**
Upon delivering the code for the first step, we will thoroughly evaluate it to ensure there are no mistakes and that it meets the requirements of achieving at least 0.4 precision and 0.1 recall as specified. Once the code passes this evaluation and meets the criteria, we will transfer 40% of the agreed payment, as this milestone is critical to proceed with the project. This milestone is critical to proceed with the project.

The remaining 60% of the payment will be made upon delivery of the final product, assuming all performance objectives have been met as agreed. Unlike the first step, which is focused on evaluating the feasibility of the project's success, the final product will need to meet additional requirements that will be discussed in the "Delivery Requirements" section. These additional demands are essential for the project's completion.

**Note on Success-Based Payment**
Our payment structure is based on success. We will only make the first payment if the initial step is completed successfully (meeting the minimum precision and recall targets and passing our evaluation). Likewise, the final payment will be made only upon successful project completion, including the additional demands outlined for the final product.

If the project is successful, we foresee an ongoing need for support and improvements to the product. In such a case, we would be happy to discuss a monthly payment arrangement for continued support, with the amount to be agreed upon by both parties. Additionally, we have more projects and would be keen to establish a long-term working relationship if this project succeeds.

**Delivery Requirements**

**First Step Delivery**

The first step of the project is focused solely on demonstrating the possibility of success. At this stage, the objective is to show that the pipeline can achieve at least 0.4 precision and 0.1 recall. No additional components or functionalities are required beyond this demonstration of the pipeline's capability.

**Final Delivery**

The final delivery must be more comprehensive and include the pipeline and the ability to save the trained model as a file at the end of the process. If any pre-processing, feature selection, or feature engineering steps are implemented, these must also be saved in a dedicated folder within the project. This ensures that all necessary steps can be reloaded and consistently used across different datasets or applications.

The final product must include a function that allows for loading the saved model and any associated pre-processing or feature selection steps. This function should accept new data and make predictions based on the previously saved model and pre-processing pipeline, making the model operational and easy to use with new inputs.

In addition, the final delivery must include comprehensive documentation of the model and the project. This should consist of model documentation that provides detailed explanations of the pipeline, pre-processing steps, model choices, and any feature selection or engineering decisions made. Additionally, a word document should be provided that explains the entire project, including how the pipeline is structured, how to use the model, and the role of each component.

The final product must be fully organized, well-documented, and properly structured. It should be easy to navigate and use, with all components clearly labeled and integrated, ensuring the pipeline is ready for deployment and further development.

**Data Privacy and Confidentiality:** The datasets provided for this challenge are private and contain sensitive information reflective of real-world events. As part of your participation in this challenge, you are required to uphold the highest standards of data privacy and confidentiality.

**Non-Disclosure:** You are strictly prohibited from sharing, copying, or distributing any part of the data to third parties. This restriction applies to both the raw data and any derivatives or analyses you produce from it.

**Data Security:** Adequate measures must be taken to ensure the security and integrity of the data throughout the duration of this project. This includes secure storage practices and the proper handling of data during analysis to prevent unauthorized access.

**Compliance:** Ensure compliance with all applicable laws and regulations regarding data protection and privacy. This responsibility includes adhering to guidelines set forth in agreements and documents signed at the outset of this project.

By participating in this challenge, you agree to abide by these conditions. Violations of these guidelines will result in immediate disqualification from the challenge, potential legal action, and forfeiture of any claims to compensation or rewards.

These considerations are essential for maintaining the trust and integrity of our data and our collaborative relationships. We appreciate your cooperation in adhering to these standards.

**Thank You**

We deeply appreciate your participation in the Data Science Pipeline Development. Your efforts and innovation are crucial to advancing our goals. Thank you for dedicating your time and expertise to this project. We look forward to seeing your creative solutions and to fostering a successful, long-term partnership.