leaders

**SCENARIO EXERCISE**

# The importance of data quality in AI-targeted marketing campaigns

# CONTENTS

Co-funded by the European Union

# Abstract

**Type of OER:**

- Scenario Exercise

**Goal/Purpose:**

- Raise awareness about the importance of data quality for implementing automated decision making algorithms (ADMS), particularly in the marketing sector.

**Expected Learning Outcomes:**

- The student will be able to implement measures to address bias in customer behaviour predictions.

**Suggested Methodological Approach (Case-Based Learning, Problem-Based Learning...):**

- Problem-based learning

**Keywords:**

- Data quality, biased algorithms, unbalanced data, marketing campaigns, ADMS.

# Introduction

# Introduction

**Context**: predicting bank marketing success using machine learning.
- The dataset used in this project is the Bank Marketing Dataset from the UCI Machine Learning Repository. It contains details about clients contacted in a marketing campaign and whether they subscribed to a term deposit.
- Dataset Description: Bank Marketing Dataset – UCI (bank.csv included in the OER materials).
- Read carefully each variable and understand its meaning.

**Key Problem Statement**
- How can we predict whether a client will subscribe to a term deposit based on their profile and past campaign interactions?
- **Challenge**: The dataset is highly imbalanced, with far fewer clients subscribing to a deposit, leading to potential bias in the model.

# Tools presentation

# Tools used in this scenario exercise

- **Python**: the most used programming language in Data Science
- **A python editor**: Google Colab or Jupyter Notebook
- **Libraries and Frameworks**
  - Pandas: data manipulation and cleaning
  - Scikit-learn: model training, evaluation, and data preprocessing
  - XGBoost: optimized gradient boosting framework
  - imbalanced-learn (SMOTE): to address class imbalance issues
  - Matplotlib/Seaborn: visualization and data insights

- **Techniques Applied**
  - Data Cleaning & Feature Engineering
  - Handling Class Imbalance (SMOTE, Undersampling, Class Weights)
  - Model Training (RandomForest, XGBoost)
  - Performance Evaluation (Recall, Precision, and F1-score)

# Hands-on activities

# Solution Description (follow the provided notebook)

## DATA PREPROCESSING

- Removed variables such as duration, campaign, and pdays that are only known after contacting clients to avoid data leakage.

- Encoded categorical variables and scaled numerical data for better model performance.

## HANDLING CLASS IMBALANCE
Applied:
- Class Weights in RandomForest to penalize errors on the minority class.
- SMOTE (Synthetic Minority Over-sampling Technique) to artificially increase samples in the minority class.
- Undersampling to reduce the size of the majority class.

## MODEL TRAINING

Tested multiple models:
- RandomForest for baseline prediction.
- XGBoost for improved performance, optimized with Early Stopping and Feature Selection to reduce training time.

## EVALUATION

Evaluated results using:
- Recall for minority class detection.
- F1-Score for balanced accuracy between precision and recall.

# Conclusion

# Conclusion

**Key Insights**
- Bias in Data matters: the imbalance in the dataset led initial models to ignore the minority class (clients subscribing to deposits).
- Balancing Techniques are key: undersampling, Class Weights, and SMOTE improved recall significantly, albeit with trade-offs in overall accuracy.
- XGBoost with Feature Selection: by reducing the number of features and adding early stopping, XGBoost improved performance without compromising efficiency.

**Key Lesson**
- Ethical AI design requires thoughtful dataset preparation, fair evaluation metrics, and awareness of potential bias in outcomes.

# References

# References

- Bank Marketing Dataset - UCI Machine Learning Repository: https://archive.ics.uci.edu/dataset/222/bank+marketing

- Scikit-learn Documentation: https://scikit-learn.org/

- XGBoost Documentation: https://xgboost.readthedocs.io/en/stable/

- Imbalanced-learn Documentation: https://imbalanced-learn.org/

# Complementary material

Click to type...

leaders

# Jupyter Notebook (IPYNB) File

- A detailed Python Notebook with commented code that walks through each step of the process is included.
- The notebook contains:
  - Data cleaning and preprocessing steps.
  - Feature engineering and variable selection logic.
  - Implementation of different balancing techniques (SMOTE, Class Weights, and Undersampling).
  - Model training with RandomForest and XGBoost.
  - Evaluation metrics and insights from the results.

"

As a technologist, I see how AI and the fourth industrial revolution will impact every aspect of people's lives.

Fei Fei Li

**leaders**

**Follow our journey**

www.aileaders-project.eu