**Budapest University of Technology and Economics**

**Faculty of Electrical Engineering and Informatics**

**Department of Artificial Intelligence and Systems Engineering**

# ML Techniques for Detecting Credit Card Fraud in Highly Imbalanced Datasets

**Cuadros Rivas, Alejandra Paola**

**KK5459**

**Field: Computer Engineering**

**Specialization: Data Science and Artificial Intelligence**

**Professor: Antal, Péter**

**2024/ Winter semester**

## I.  Abstract

The detection of fraudulent transactions in credit card data represents a critical challenge in financial security, demanding analytical solutions to protect consumers. This project focuses on developing a predictive model capable of identifying fraudulent activities within a dataset characterized by a severe class imbalance, where fraudulent transactions constitute only 0.172% of the data. The dataset, encompassing transactions from European cardholders in September 2013, includes 284,807 transactions with 492 frauds, and 30 features primarily derived through PCA transformation, alongside 'Time' and 'Amount' which have not been transformed.

We employed various machine learning techniques to address the imbalance and improve fraud detection. Preprocessing steps included scaling the Amount feature, adding Time_of_day to capture temporal patterns, and applying SMOTE, under-sampling, and a combination of both to balance the dataset. Among the models tested, Gradient Boosting with SMOTE emerged as the most effective, achieving high precision (0.99), recall (0.99), and ROC-AUC (0.9998). Dimensionality reduction through PCA reduced the dataset's features from 30 to 17 while maintaining strong performance, highlighting its potential for simplifying complex datasets. Explainability tools like SHAP and LIME played a critical role in identifying key features (V14, V4, and V12) and providing interpretable predictions for fraud and non-fraud transactions, fostering trust in the models.

Additionally, a real-time simulation pipeline was implemented to mimic real-world conditions, demonstrating the model's ability to handle live transaction streams with minimal latency. The findings underscore the importance of data balancing techniques, robust supervised models like Gradient Boosting, and explainability tools in building effective fraud detection systems.

Future work includes exploring dynamic cost-sensitive learning, deploying real-time adaptive models, and scaling the solution to larger, more diverse datasets to enhance generalizability and robustness.