

# Hackathon BI Pipeline – Data Science Project

## Watsonx – Banking Fraud Detection

### Team 20

- Lisa NACCACHE (DIA 4) – Team Leader
- Hiba NEJJARI (DIA 4)
- Neil MAHCER (DIA 4)
- Wendy DUONG (DIA 4)
- Cyprien MOUTON (DIA 4)
- Safa HORMI BOUAICHI (DIA 3)

### 1. Introduction and Context

With the rapid growth of digital banking transactions, financial institutions are increasingly exposed to fraud attempts. Fraudsters continuously adapt and develop more sophisticated strategies, making fixed-rule detection systems less effective.

For this hackathon, IBM provided a real-world dataset derived from banking transactions recorded between 2016 and 2018, enabling us to develop, train, and evaluate fraud detection models.

Our objective is to build an Artificial Intelligence model capable of predicting whether a transaction is fraudulent (1) or legitimate (0), in order to support banking security teams in real-time decision-making.

### 2. Problem Statement

Credit card fraud leads to significant financial losses for banks and their customers. Detection must be fast to immediately block suspicious transactions, reliable to avoid false alerts that inconvenience customers, and capable of generalizing since fraudsters continuously change their behavior.

Our problem is therefore: How can we develop a Machine Learning model capable of effectively identifying fraudulent transactions?

### 3. Description of the IBM Dataset

The datasets we have at our disposal are the following:

- **transactions\_train.csv**, which contains all the banking transactions (amount, date, location, card used, etc.). This will be our main dataset that will allow us to detect unusual behaviors, such as abnormal amounts or transactions in different cities.
- **train\_fraud\_labels.json**, which indicates whether each transaction in the training set is fraudulent (1) or not (0). This file provides the ground truth for training the model.
- **cards\_data.csv**, which contains technical information about the cards. This helps us understand whether a card is at risk. For example, a card reported on the dark web is a strong indicator of fraud.
- **user\_data.csv**, which provides customer profile information (age, income, location, credit score, number of cards). This allows us to compare a transaction against the customer's typical behavior — for instance, low income + large purchase = suspicious.
- **mcc\_codes.json**, which provides a classification of merchants (restaurants, online gaming, transportation, etc.). This allows us to analyze spending habits and identify purchases in high-risk sectors (e.g., casinos have a higher fraud rate).
- **evaluation\_features.csv**, which contains unlabeled transactions that must be predicted for submission. This file will not be used for training but only for generating the final prediction file.

### 4. Cold Start Challenge

Unlike many classification problems, the customers present in the training set are not the same as those in the evaluation set.

This means that our model must:

- Generalize without memorizing the specific behaviors of certain users.
- Rely more on transaction patterns rather than user identities.
- Be resilient to new or previously unseen behaviors.

## 5. Methodology

As a group of six students from the Data & Artificial Intelligence major, we chose to organize ourselves as a real project team by assigning roles according to each member's skills, ensuring complementarity between the technical, analytical, and visual aspects. In parallel, we set up a [GitHub](#) repository to centralize and coordinate our work.

### Business Understanding

- Identify the main challenges related to banking fraud and the constraints of the project (cold start, class imbalance).
- Define the model's objectives and evaluation metrics (AUROC, AUPRC, F1-score).

### Data Understanding

- Study the structure of the dataset, detect anomalies and missing values, and understand typical transactional behaviors.

### Data Preparation

- Clean, merge, and create new features (time-based variables, categorical encodings, etc.) while preventing any data leakage.

### Modeling

- Train a Machine Learning model (LightGBM) capable of handling strong class imbalance and generalizing to new customers.
- Implement temporal validation: train on 2016–2017 data and validate on 2018 data.

### Evaluation & Generalization

- Measure performance on unseen customers and transactions to ensure robustness.
- Adjust the classification threshold according to the most relevant metric (F1 / AUPRC).

### Visualization & Deployment

- Create an interactive dashboard to present the results and key model indicators.
- Generate the final submission file in the required format (transaction\_id, fraud\_prediction).
- Write a clear synthesis for the final report.

All these tasks are documented in detail in our GitHub repository, in the folder:  
Management > [tasks.md](#)

## Team Task Distribution

Phase	Main tasks	Responsible members
<b>1</b> Introduction & Context	Presentation of the project and the IBM dataset	Lisa & Cyprien
<b>2</b> Data Exploration	Exploratory analysis and initial visualization	Hiba
<b>3</b> Data Preparation	Cleaning, merging, feature engineering	Hiba & Lisa & Safa
<b>4</b> Machine Learning Modeling	Training, model selection, tuning	Neil & Wendy
<b>5</b> Evaluation & Generalization	Robustness testing and comparison	Neil & Wendy
<b>6</b> Data Visualization & Dashboard	Creation of the dashboard in Power BI	Safa
<b>7</b> Generation of the submission file	Production of the final CSV for evaluation	Hiba & Cyprien
<b>8</b> Documentation & Presentation	Report and oral presentation	All members