

# Data Dictionary – Flat User Credit Scoring Dataset (Labeled)

This document describes the final flat, labeled dataset used for training and evaluating AI-based Credit Scoring models in the Vietnamese fintech context.

## 1. Dataset Overview

Dataset name: flat\_user\_credit\_scoring\_labeled.csv

Grain: One row per user

Purpose: Model training and evaluation for credit risk classification.

## 2. User Identifier and Demographics

- user\_id: Anonymous user identifier, randomly generated.
- user\_age: Age of the user, ranging from 18 to 65.
- user\_region: Region of residence (HN, HCM, DN, CT, OTHER).

## 3. Telco Aggregated Features

- telco\_account\_age\_days: Number of days the telecom account has been active.
- telco\_avg\_revenue\_mean: Average monthly telecom spending.
- telco\_avg\_revenue\_std: Volatility of telecom spending.
- telco\_recharge\_count\_mean: Average monthly recharge count.
- telco\_recharge\_amount\_mean: Average monthly recharge amount.
- telco\_mobile\_data\_mb\_mean: Average monthly mobile data usage.
- has\_telco\_data: Flag indicating availability of telco data.

## 4. Academic Features

- edu\_highest\_level: Highest education level achieved.
- edu\_gpa\_band: GPA category band.
- edu\_graduation\_status: Graduation status.
- edu\_institution\_tier: Quality tier of the institution.
- edu\_major\_group: Broad field of study.
- has\_academic\_data: Flag indicating availability of academic data.

## 5. E-Wallet Aggregated Features

- wallet\_txn\_count: Total number of wallet transactions.
- wallet\_avg\_amount: Average transaction amount.
- wallet\_max\_amount: Maximum transaction amount.

- `wallet_total_amount`: Total transaction amount.
- `wallet_large_txn_ratio`: Ratio of large transactions above 2 million VND.
- `wallet_failure_rate`: Ratio of failed transactions.
- `has_ewallet_data`: Flag indicating availability of e-wallet data.

## 6. Target Label Definition

`target`: Binary credit risk label.

0 = Good customer (low risk).

1 = Bad customer (high risk).

The label is synthetically generated using stability, capacity, and behavioral risk indicators. No future information is used, and the bad rate is approximately 35 percent.