

**КИЇВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ ТАРАСА  
ШЕВЧЕНКА**

**Факультет інформаційних технологій**

Кафедра технологій управління

Спеціальність 122 – Комп'ютерні науки,  
освітня програма «Інформаційна аналітика та впливи»

**КВАЛІФІКАЦІЙНА РОБОТА МАГІСТРА**

на тему:

**«Інтелектуальні системи виявлення шахрайських операцій в онлайн-  
платежах на підґрунті моделей машинного навчання»**

**Студента 2-го курсу групи ІАВ-21**

**Науковий керівник:**

Супруненко Максим Ігорович

к.е.н, доцент

Мірошниченко І.В

(прізвище, ім'я, по батькові)

(науковий ступінь, вчене звання,  
прізвище, ім'я, по батькові)

\_\_\_\_\_

\_\_\_\_\_

(підпис студента) (дата)

\_\_\_\_\_

\_\_\_\_\_

(підпис) (дата)

**Попередній захист:**

---

(Висновок: «До захисту в Екзаменаційній комісії»)

Завідувач кафедри  
технологій управління

\_\_\_\_\_

\_\_\_\_\_

\_\_\_\_\_

(підпис) (прізвище, ініціали) (дата)

**Київ – 2026**

**КИЇВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ ТАРАСА  
ШЕВЧЕНКА**

**Факультет інформаційних технологій**

Кафедра технологій управління  
Освітній рівень Магістр  
Спеціальність 122 - Комп'ютерні науки  
Освітньо-наукова програма Інформаційна аналітика та впливи

**ЗАТВЕРДЖУЮ**

Завідувач кафедри

«\_\_\_» \_\_\_\_\_ 2026 року

**ЗАВДАННЯ  
НА ВИКОНАННЯ КВАЛІФІКАЦІЙНОЇ РОБОТИ**

Студент Супруненко Максим Ігорович Група ІАВ-21

**1. Тема кваліфікаційної роботи**

Інтелектуальні системи виявлення шахрайських операцій в онлайн-платежах на підґрунті моделей машинного навчання

Затверджена наказом по від «\_\_\_» \_\_\_\_\_ 2026 р. № \_\_\_\_.

**2. Строк подання студентом готової роботи – «\_\_\_» \_\_\_\_\_ 2026 р.**

**3. Цільова установка та вихідні дані до роботи**

Розробка та дослідження комплексної системи виявлення шахрайських акаунтів у платіжному сервісі на основі методів машинного навчання (логістична регресія, градієнтний бустінг LightGBM та XGBoost, нейронні мережі MLP, граф-нейронні мережі GraphSAGE та GAT, стакінг-ансамблі з мета-навчанням) з метою максимізації F1-score за умови сильного дисбалансу класів (співвідношення 25:1). Вхідні дані: анонімізований датасет SKELAR × mono AI Competition 2025 — 564 830 користувачів, 4,5 млн транзакцій, 283 сконструйовані ознаки (базові агрегати, velocity-ознаки, error-ознаки, графові

LOO-ознаки, target encoding) та card-sharing граф із 1,8 млн ребер, побудований на основі спільних ідентифікаторів карток.

#### **4. Зміст роботи**

Робота містить вступ, огляд теоретичних засад виявлення фінансового шахрайства з аналізом сучасних методів машинного навчання та граф-нейронних мереж, а також постановку завдання. Описано дані змагання, проведено розвідувальний аналіз і представлено методологію конструювання ознак та побудови card-sharing графу. Наведено математичне обґрунтування ключових методів і запропоноване покращення шляхом інтеграції GNN у стакінг-ансамбль з мета-навчанням. Після навчання моделей проведено аналіз отриманих результатів. Представлено стратегію бізнес-інтеграції з трирівневою системою прийняття рішень та обґрунтуванням економічного ефекту. У заключенні подані висновки, а також список використаних джерел та додатки.

#### **5. Перелік графічного матеріалу (слайдів)**

Архітектура комплексної системи виявлення шахрайства (стакінг-ансамбль із 19 моделей), прогресія F1-score протягом дослідження (від 0,5960 до 0,8877), порівняння запропонованого підходу з альтернативними методами (Rule-based, Isolation Forest, Label Propagation, LightGBM, GraphSAGE окремо), сегментний аналіз card-sharing графу та fraud rate за сегментами (зв'язані vs ізольовані вузли), аналіз важливості ознак за категоріями (feature importance), трирівнева стратегія впровадження антифрод-системи у платіжний сервіс (зелена / жовта / червона зони).

## 6. Календарний план виконання роботи:

№ п/п	Назва частин роботи	%	Виконання роботи	
			За планом	Фактично
1.	Вибір теми дипломної роботи	3	01.10.26	01.10.26
2.	Протокол кафедри ТУ про затвердження тем дипломних робіт та призначення наукових керівників	2	27.12.26	27.12.26
3.	Формування переліку Нормативних матеріалів, літератури з проблематики дипломної роботи	10	08.01.26	07.01.26
4.	Складання розгорнутого плану кваліфікаційної роботи	5	18.01.26	18.01.26
5.	Ознайомлення наукового керівника з розгорнутим планом кваліфікаційної роботи. Внесення змін.	5	19.01.26 - 20.01.26	20.11.26
6.	Підготовка розділу 1	10	12.02.26	13.02.26
7.	Підготовка розділу 2	14	08.03.26	08.03.26
8.	Підготовка розділу 3	14	01.04.26	01.04.26
9.	Підготовка розділу 4	13	20.04.26	20.04.26
10.	Оформлення кваліфікаційної роботи. Підготовка висновків і пропозицій	15	30.04.26	30.04.26
11.	Передача кваліфікаційної роботи науковому керівникові	2	04.05.26	04.05.26
12.	Передача кваліфікаційної роботи рецензенту для рецензування	2	07.05.26	11.05.26
13.	Попередній захист кваліфікаційної роботи	5	11.05.26	17.05.26

Дата видачі завдання «\_\_» \_\_\_\_\_ 20\_\_ р.

Керівник роботи к.е.н., доцент кафедри технологій управління МІРОШНИЧЕНКО І.В.

\_\_\_\_\_  
(підпис)

Завдання прийняв до виконання студент групи ІАВ-21 Супруненко Максим Ігорович

\_\_\_\_\_  
(підпи

## ЗМІСТ

АНОТАЦІЯ.....	7
ПЕРЕЛІК ВИКОРИСТАНИХ СКОРОЧЕНЬ .....	9
ВСТУП.....	10
РОЗДІЛ 1. ТЕОРЕТИЧНІ ЗАСАДИ ВИЯВЛЕННЯ ФІНАНСОВОГО ШАХРАЙСТВА .....	14
1.1. Характеристика проблеми виявлення шахрайських операцій в онлайн- платежах та аналіз предметної області .....	14
1.2. Методи машинного навчання для класифікації шахрайських акаунтів ....	24
1.3. Ансамблеві методи та градієнтний бустинг у виявленні шахрайства .....	26
1.4. Застосування нейронних мереж для табличних даних у антифроді.....	28
Висновки до розділу 1 .....	29
РОЗДІЛ 2. ДАНІ ТА МЕТОДОЛОГІЯ ДОСЛІДЖЕННЯ .....	31
2.1. Опис змагання та джерел даних дослідження виявлення шахрайства .....	31
2.2. Структура та характеристики даних платіжного сервісу .....	32
2.3. Розвідувальний аналіз даних для виявлення патернів шахрайської поведінки .....	34
2.4. Конструювання ознак для класифікації шахрайських акаунтів.....	53
2.5. Побудова card-sharing графу для моделювання мережових зв'язків між акаунтами.....	57
2.6. Граф-нейронні мережі та архітектура GraphSAGE для виявлення мережевого шахрайства .....	59
2.7. Методологія оцінювання моделей та запобігання витоку даних.....	61
2.8. Focal Loss як метод роботи з дисбалансом класів у задачах виявлення шахрайства.....	64
2.9. Стакінг та мета-навчання для побудови ансамблю моделей .....	66

2.10. Організація командної роботи та інструменти розробки антифрод-системи.....	68
Висновки до розділу 2.....	69
<b>РОЗДІЛ 3. МОДЕЛІ, ЕКСПЕРИМЕНТАЛЬНІ РЕЗУЛЬТАТИ ТА СТРАТЕГІЯ ВПРОВАДЖЕННЯ</b> .....	<b>71</b>
3.1. Базові моделі класифікації шахрайських акаунтів .....	71
3.2. Застосування LightGBM з Focal Loss для роботи з дисбалансом класів ...	73
3.3. Мульти-модельний пайплайн навчання та селекції базових моделей.....	75
3.4. Навчання GraphSAGE на card-sharing мережі користувачів .....	77
3.5. Стакінг-ансамбль з forward selection і мета-моделлю XGBoost.....	79
3.6. Аналіз помилок класифікації та виявлення «тихого фроду» .....	101
3.7. Трирівнева стратегія впровадження антифрод-системи .....	105
3.8. Моніторинг concept drift та стратегія перенавчання моделі.....	107
3.9. Інтерпретація моделі та аналіз важливості ознак .....	107
3.10. Економічний ефект та план впровадження антифрод-системи .....	111
Висновки до розділу 3 .....	112
<b>ВИСНОВКИ</b> .....	<b>113</b>
<b>СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ</b> .....	<b>116</b>
<b>ДОДАТКИ</b> .....	<b>119</b>
Додаток А. Конфігурації базових моделей .....	119
Додаток Б. Архітектура GraphSAGE .....	120
Додаток В. Оптимальний ансамбль.....	121
Додаток Г. Детальні результати крос-валідації .....	122
Додаток Д. Кореляційна матриця OOF-предиктів.....	124
Додаток Е. Feature Engineering: повний список ознак	<b>Error! Bookmark not defined.</b>
Додаток Ж. Програмний код ключових компонентів	<b>Error! Bookmark not defined.</b>

## АНОТАЦІЯ

### КИЇВСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ ТАРАСА ШЕВЧЕНКА

#### Факультет інформаційних технологій

Кафедра технологій управління

Спеціальність 122 - Комп'ютерні науки,

освітня програма "Інформаційна аналітика та впливи"

Кваліфікаційна робота магістра Супруненко Максима Ігоровича.

Тема роботи – «Інтелектуальні системи виявлення шахрайських операцій в онлайн-платежах на підґрунті моделей машинного навчання».

Мета кваліфікаційної роботи магістра – розробка та дослідження інтелектуальних систем виявлення шахрайських акаунтів у платіжному сервісі на основі методів машинного навчання, включаючи табличні моделі градієнтного бустінгу, граф-нейронні мережі та стакінг-ансамблі, з досягненням максимального значення F1-score на реальних даних змагання SKELAR × mono AI Competition.

Об'єкт дослідження – процес виявлення шахрайських акаунтів у платіжному сервісі з використанням методів машинного навчання.

Предмет дослідження – методи та алгоритми машинного навчання для бінарної класифікації за ознакою шахрайської діяльності, включаючи градієнтний бустінг, граф-нейронні мережі, ансамблеві методи та підходи до конструювання ознак.

Наукова новизна роботи – запропоновано комплексний підхід, що поєднує табличні моделі градієнтного бустінгу з граф-нейронними мережами GraphSAGE на card-sharing графі та мульти-модельний стакінг з forward selection, що дозволило досягти OOF F1-score 0.8877 на реальних даних.

У роботі досліджуються сучасні підходи до виявлення фінансового шахрайства. Розроблено систему конструювання 280+ ознак, побудовано card-

sharing граф із 564 830 вузлів та 1,8 млн ребер, реалізовано стакінг-ансамбль із 19 моделей з rank-нормалізацією. Запропоновано трирівневу стратегію бізнес-інтеграції антифрод-моделі.

Кваліфікаційна робота складається зі вступу, основної частини, яка включає три розділи, висновків та списку використаних джерел. Робота містить рисунки, таблиці та додатки.

*Ключові слова: машинне навчання, виявлення шахрайства, граф-нейронні мережі, GraphSAGE, градієнтний бустінг, LightGBM, стакінг, card-sharing граф, F1-score, Focal Loss, антифрод.*

## ПЕРЕЛІК ВИКОРИСТАНИХ СКОРОЧЕНЬ

- Програмний інтерфейс застосунку (далі – API, Application Programming Interface)
- Площа під ROC-кривою (далі – AUC, Area Under ROC Curve)
- Бінарна перехресна ентропія (далі – BCE, Binary Cross-Entropy)
- Перехресна валідація (далі – CV, Cross-Validation)
- Глибинне навчання (далі – DL, Deep Learning)
- F1-міра — гармонічне середнє точності та повноти (далі – F1-score)
- Функція втрат з фокусуванням на складних прикладах (далі – FL, Focal Loss)
- Градiєнтний бустiнг (далі – GBM, Gradient Boosting Machine)
- Загальний регламент захисту даних ЄС (далі – GDPR, General Data Protection Regulation)
- Граф-нейронна мережа (далі – GNN, Graph Neural Network)
- Графічний процесор (далі – GPU, Graphics Processing Unit)
- Архітектура GNN з вибірковою агрегацією сусідів (далі – GraphSAGE, Graph Sample and Aggregate)
- Кваліфікаційна робота магістра (далі – KPM)
- Бібліотека градієнтного бустінгу від Microsoft (далі – LightGBM)
- Машинне навчання (далі – ML, Machine Learning)
- Багатошаровий перцептрон (далі – MLP, Multi-Layer Perceptron)
- Передбачення на валідаційних фолдах (далі – OOF, Out-of-Fold)
- Метод головних компонент (далі – PCA, Principal Component Analysis)
- Крива характеристик приймача (далі – ROC, Receiver Operating Characteristic)
- Елементи матриці плутанини (далі – TP/FP/FN/TN, True/False Positive/Negative)
- Відеопам'ять (далі – VRAM, Video Random Access Memory)
- Бібліотека градієнтного бустінгу (далі – XGBoost, eXtreme Gradient Boosting)

## ВСТУП

### **Актуальність теми дослідження.**

Стрімкий розвиток цифрових платіжних технологій у XXI столітті створив безпрецедентні можливості для фінансового обслуговування населення. Мобільні платіжні сервіси, такі як Apple Pay, Google Pay щорічно обробляють мільярди транзакцій, формуючи нову парадигму фінансових відносин. Водночас зростання обсягів електронних платежів супроводжується пропорційним збільшенням масштабів фінансового шахрайства. За даними Nilson Report, глобальні втрати від карткового шахрайства у 2023 році перевищили 33 мільярди доларів США, демонструючи щорічне зростання на 12-15%. Для платіжних сервісів, які обслуговують мільйони клієнтів, навіть незначне підвищення точності виявлення шахрайських акаунтів здатне зберегти десятки мільйонів гривень і підвищити довіру користувачів до цифрових фінансових послуг.

Традиційні підходи до виявлення шахрайства, засновані на експертних правилах та порогових значеннях, виявляються недостатньо ефективними в умовах постійної адаптації зловмисників. Сучасні шахрайські схеми використовують складні мережеві структури, де група зловмисників ділиться платіжними картками та використовує спільні реквізити для маскування своєї діяльності. Класичні табличні моделі машинного навчання, хоча й демонструють значний прогрес порівняно з правилами, не здатні повноцінно моделювати такі мережеві взаємодії. Саме тому актуальним є дослідження граф-нейронних мереж (GNN) як інструменту для виявлення фроду на рівні мережевих зв'язків між користувачами.

Дана кваліфікаційна робота магістра виконана в рамках змагання SKELAR × mono AI Competition 2025 — реального індустріального хакатону з виявлення шахрайських акаунтів у мобільному платіжному сервісі monobank.

### **Зв'язок роботи з науковими програмами, планами, темами.**

Кваліфікаційна робота магістра виконана на кафедрі технологій управління факультету інформаційних технологій Київського національного університету імені Тараса Шевченка в рамках кафедральної теми 0121U107799 «Розробка інформаційно-аналітичних інструментів управління портфелями проєктів і програм в інтегрованих функціональних середовищах».

### **Мета і завдання дослідження.**

Метою кваліфікаційної роботи магістра є розробка та дослідження комплексної системи виявлення шахрайських акаунтів у мобільному платіжному сервісі на основі методів машинного навчання, включаючи табличні моделі градієнтного бустінгу, граф-нейронні мережі та стакінг-ансамблі, з досягненням максимального значення F1-score.

Для досягнення поставленої мети необхідно вирішити наступні завдання:

- 1) дослідити сучасні підходи до виявлення шахрайства в онлайн-платежах та обґрунтувати вибір комплексної архітектури системи;
- 2) виконати розвідувальний аналіз даних і сконструювати інформативні ознаки для класифікації шахрайських акаунтів;
- 3) побудувати card-sharing граф та навчити граф-нейронні мережі GraphSAGE і GAT для класифікації вузлів;
- 4) розробити стакінг-ансамбль базових моделей із мета-моделлю XGBoost та rank-нормалізацією прогнозів;
- 5) виконати валідацію моделі, проаналізувати помилки класифікації та розробити стратегію її інтеграції у платіжний сервіс.

### **Об'єкт дослідження.**

Об'єктом дослідження є процес виявлення шахрайських акаунтів у мобільному платіжному сервісі з використанням методів машинного навчання.

### **Предмет дослідження.**

Предметом дослідження є методи та алгоритми машинного навчання для бінарної класифікації користувачів мобільного платіжного сервісу за ознакою шахрайської діяльності.

### **Методи дослідження.**

У роботі використано комплекс методів: градієнтний бустінг (LightGBM, XGBoost) з різними функціями втрат, у тому числі Focal Loss; граф-нейронні мережі (GraphSAGE); нейронні мережі (MLP); стакінг з логістичною регресією як мета-модель; out-of-fold (OOF) передбачення; forward selection для оптимізації складу ансамблю; статистичний аналіз.

### **Наукова новизна одержаних результатів.**

Наукова новизна одержаних результатів полягає в тому, що в роботі вперше запропоновано комплексний підхід до виявлення шахрайських акаунтів, що поєднує табличні моделі градієнтного бустингу (LightGBM, XGBoost) з граф-нейронними мережами GraphSAGE та GAT на card-sharing графі через стакінг-ансамбль з XGBoost-мета-моделлю та rank-нормалізацією OOF-прогнозів. Окремими елементами наукової новизни є: методологія побудови card-sharing графу (564 830 вузлів, ~1,8 млн ребер) на основі спільних карткових ідентифікаторів та реквізитів карткотримача; сегментовані пороги класифікації для зв'язаних і ізольованих сегментів графу; адаптація Focal Loss до задачі класифікації акаунтів з дисбалансом класів 25:1.

### **Практичне значення одержаних результатів.**

Розроблена модель досягає OOF F1-score 0.8877, що відповідає recall ~86% — система ідентифікує 86% шахрайських акаунтів. Запропонована трирівнева стратегія бізнес-інтеграції дозволяє поетапно впровадити систему.

### **Апробація результатів.**

Результати апробовано у рамках змагання SKELAR × mono AI Competition 2026 на реальних даних мобільного платіжного сервісу monobank.

**Публікації.** За темою дисертаційного дослідження та суміжними тематиками опубліковано тези двох доповідей у співавторстві з науковим керівником к.е.н., доцентом Мирошніченком І. М. на факультетських науково-практичних конференціях Київського національного університету імені Тараса Шевченка: «Photorealistic Synthetic Data Generation for Robust Ukrainian Fiscal Receipt OCR» та «Information Analytics for Forming and Supporting Productive Habits». Окрім того, результати дисертаційного дослідження пройшли практичну апробацію у рамках хакатону SKELAR × mono AI Competition 2026 на реальних анонімізованих даних мобільного платіжного сервісу.

**Структура та обсяг роботи.**

Кваліфікаційна робота магістра викладена на 116 сторінках, складається зі вступу, трьох розділів, висновків, списку використаних джерел із 38 найменувань та 8 додатків. Робота містить 26 рисунків, 26 таблиць та 200 формул.

# РОЗДІЛ 1. ТЕОРЕТИЧНІ ЗАСАДИ ВИЯВЛЕННЯ ФІНАНСОВОГО ШАХРАЙСТВА

## 1.1. Характеристика проблеми виявлення шахрайських операцій в онлайн-платежах та аналіз предметної області

### Загальний огляд ландшафту цифрового шахрайства

Глобальний ландшафт цифрового шахрайства зазнав кардинальних змін протягом останнього десятиліття. Якщо у 2010-х роках основним вектором атак було шахрайство з фізичними платіжними картками (card-present fraud), то з масовим переходом на безконтактні та мобільні платежі акцент змістився на card-not-present (CNP) шахрайство, яке у 2023 році становило понад 70% усіх карткових втрат. Мобільні платіжні сервіси стали привабливою ціллю для зловмисників з кількох причин: 1) спрощена процедура реєстрації порівняно з традиційними банками; 2) можливість швидкого масштабування операцій через автоматизацію (боти); 3) глобальний характер сервісів, що ускладнює юрисдикційне переслідування.

За даними Juniper Research [19], глобальні втрати від онлайн-шахрайства досягнуть 48 мільярдів доларів до 2027 року. У України, де банки обслуговує понад 20 мільйонів клієнтів, навіть 0.1% рівень шахрайства означає десятки тисяч постраждалих акаунтів та мільйонні збитки. Це робить розробку ефективних антифрод-систем стратегічним пріоритетом для фінтех-компаній.

Еволюція шахрайських схем характеризується зростаючою складністю та організованістю. Сучасні шахрайські угруповання використовують: автоматизовані інструменти для масової реєстрації акаунтів (account farming); ботнети для розподілу активності між тисячами IP-адрес; соціальну інженерію для компрометації даних легітимних користувачів; криптовалюти для відмивання

отриманих коштів; dark web маркетплейси для торгівлі скомпрометованими даними карток.

У відповідь на зростаючу загрозу фінансові установи інвестують значні кошти в антифрод-технології. За даними LexisNexis Risk Solutions, середня фінтех-компанія витрачає 6-8% виручки на боротьбу з шахрайством, включаючи технології, персонал та компенсації. Методи машинного навчання стали стандартом індустрії, замінивши або доповнивши традиційні правилоподібні системи (rule-based systems). Однак «гонка озброєнь» між захисниками та атакуючими вимагає постійного вдосконалення моделей та підходів.

### **Типологія шахрайства в мобільних платіжних системах**

Шахрайство у мобільних платіжних системах має специфічну типологію, що відрізняється від традиційного банківського шахрайства. Основні типи включають:

Account farming (фермерство акаунтів) — масова реєстрація акаунтів з використанням синтетичних або вкрадених персональних даних. Метою є створення «чистих» акаунтів, які потім використовуються для шахрайських операцій або продаються на dark web. Характерні ознаки: однотипні патерни реєстрації (однакове джерело трафіку, схожі часові інтервали), відсутність або мінімальна легітимна активність після реєстрації, використання VPN/проксі для маскуванню геолокації.

Card testing (тестування карток) — використання мобільного платіжного сервісу для перевірки валідності вкрадених даних платіжних карток. Зловмисник здійснює серію дрібних транзакцій (зазвичай \$1-3) для перевірки, чи картка активна та не заблокована. Характерні ознаки: велика кількість невдалих транзакцій, малі суми, використання різних карток з одного акаунту, висока швидкість між транзакціями (автоматизація).

Friendly fraud (дружнє шахрайство) — ситуація, коли легітимний власник картки здійснює транзакцію, а потім заперечує її (chargeback), стверджуючи, що

транзакція була несанкціонованою. Цей тип особливо складний для виявлення, оскільки акаунт належить реальній людині з легітимною історією. У SKELAR × моно цей тип, ймовірно, менш поширений, оскільки дані фокусуються на акаунтному, а не транзакційному фроді.

Subscription abuse (зловживання підписками) — використання вкрадених карток для оформлення платних підписок з подальшим отриманням вигоди (перепродаж акаунтів, отримання промо-бонусів тощо). Характерні ознаки: транзакції на характерні «підписочні» суми (\$3, \$6, \$11, \$21, \$51), що відповідає популярним сервісам; повторювані транзакції з фіксованою періодичністю; високий відсоток card\_recurring транзакцій.

### **Метрики оцінювання для задач з дисбалансом**

Вибір правильної метрики оцінювання є критичним для задач із суттєвим дисбалансом класів. Стандартна accuracy (частка правильних передбачень) є оманливою: модель, що завжди передбачає негативний клас, досягне accuracy = 96.22% при fraud rate 3.78%. Тому для оцінки антифрод-моделей використовуються спеціалізовані метрики.

Precision (точність) =  $TP / (TP + FP)$  — частка справжніх шахраїв серед усіх акаунтів, які модель класифікувала як шахрайські. Висока precision означає мало хибних спрацьовувань. У бізнесу: precision = 0.87 означає, що 87% заблокованих акаунтів — справді шахраї, а 18% — легітимні користувачі, яких помилково заблокували.

Recall (повнота) =  $TP / (TP + FN)$  — частка виявлених шахраїв серед усіх справжніх шахраїв. Високий recall означає мало пропущених шахраїв. У бізнесу: recall = 0.90 означає, що модель виявляє 90% шахрайських акаунтів, пропускаючи 12%.

$$F_1 - score = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall} \quad (1.1)$$

гармонічне середнє precision та recall. F1 штрафує за дисбаланс між precision та recall: модель з precision=0.99 та recall=0.01 матиме F1=0.02, тоді як модель з precision=0.50 та recall=0.50 — F1=0.50. У змаганні SKELAR × mono F1 є основною метрикою (60% ваги).

ROC-AUC (Area Under ROC Curve) — площа під кривою, що графічно відображає залежність True Positive Rate від False Positive Rate при різних порогах класифікації. AUC  $\in [0, 1]$ , де 0.5 — випадковий класифікатор, 1.0 — ідеальний. AUC є threshold-independent метрикою, тобто оцінює якість ранжування, а не конкретного бінарного рішення. У дослідженні AUC використовується для моніторингу навчання GNN (ранній зупин по AUC).

PR-AUC (Precision-Recall AUC) — площа під кривою Precision-Recall, яка є більш інформативною за ROC-AUC для задач із сильним дисбалансом. Baseline PR-AUC = fraud\_rate = 0.0378 (випадковий класифікатор), що робить цю метрику чутливішою до реального покращення якості.

### Детальний огляд velocity-ознак

Velocity-ознаки (ознаки швидкості) є однією з найінформативніших категорій для виявлення шахрайства, оскільки вони захоплюють темпоральні патерни поведінки, що суттєво відрізняються між легітимними користувачами та шахраями. Концептуально, velocity-ознаки вимірюють «як швидко» та «наскільки інтенсивно» користувач здійснює транзакції.

Інтервали між транзакціями (time deltas). Для кожного користувача обчислюється послідовність часових інтервалів між сусідніми транзакціями:

$$\Delta t_i = t_{i+1} - t_i \quad (1.2)$$

Статистики цієї послідовності — mean, std, min, max, median, percentiles — є потужними ознаками. Шахраї характеризуються: дуже малими мінімальними інтервалами (< 10 секунд), що вказує на автоматизацію; низькою варіативністю інтервалів (std  $\rightarrow$  0), що вказує на скриптову поведінку; біфуркацією розподілу (або дуже швидко, або довгі паузи між «сесіями»).

Burst-аналіз (аналіз спалахів активності). Burst визначається як група послідовних транзакцій з інтервалом менше 60 секунд. Ознаки:  $n_{bursts}$  —

кількість burst-епізодів; `max_burst_length` — максимальна кількість транзакцій у одному burst; `avg_burst_length` — середня довжина burst; `total_burst_tx` — загальна кількість транзакцій у burst-режимі; `burst_ratio` — частка транзакцій у burst-режимі від загальної кількості. Шахраї мають значно більше burst-епізодів та довші burst-послідовності.

Нічна активність (`night_ratio`). Частка транзакцій, здійснених між 00:00 та 06:00 за UTC. Легітимні користувачі рідко здійснюють транзакції вночі (`night_ratio < 0.05`), тоді як автоматизовані шахрайські системи працюють цілодобово (`night_ratio > 0.15`). Ця ознака є четвертою за важливістю у LightGBM моделі.

Acceleration (прискорення активності). Вимірює зміну інтенсивності транзакцій з часом. Обчислюється як різниця між середньою кількістю транзакцій у другій та першій половині активного періоду. Позитивне прискорення вказує на наростання активності — характерний патерн для шахраїв, які починають з малої кількості тестових транзакцій та поступово збільшують масштаб.

## Математичні основи GraphSAGE

Крок 1 (Семплювання): Для кожного цільового вузла  $v$  відбираємо фіксовану кількість  $K^l$  сусідів з рівномірним розподілом:  $N_S^l(v) \subset N(v)$ ,  $|N_S^l(v)| = K^l$ . У нашій конфігурації  $K^1=15$ ,  $K^2=10$ .

$$a_v^l = MEAN(\{h_u^{l-1} : u \in N_S^l(v)\}) \quad (1.3)$$

Крок 3 (Оновлення):

$$h_v^l = \sigma(W^l \cdot CONCAT(h_v^{l-1}, a_v^l) + b^l) \quad (1.4)$$

Фінальне передбачення:

$$\hat{y}_v = \sigma(W_{head} \cdot h_v^L + b_{head}) \quad (1.5)$$

де  $W_{head} \in R^{1 \times d-L}$ ,  $\sigma$  — сигмоїдна функція. Функція втрат — зважена Binary Cross-Entropy:  $L = -\sum_v [w_{pos} \cdot y_v \cdot \log(\hat{y}_v) + (1-y_v) \cdot \log(1-\hat{y}_v)]$ , де  $w_{pos} = 25.5$  — вага позитивного класу.

$$L = -\sum_v [w_{pos} \cdot y_v \cdot \log(\hat{y}_v) + (1 - y_v) \cdot \log(1 - \hat{y}_v)] \quad (1.6)$$

### Порівняння методів боротьби з дисбалансом

Для системного порівняння методів боротьби з дисбалансом класів розглянемо їх у контексті нашої задачі з fraud rate 3.78%.

Class weighting (зважування класів).

$$scale\_pos\_weight = \frac{N_{negative}}{N_{positive}} \approx 25 \quad (1.7)$$

кожен позитивний приклад вагується у 25 разів більше при обчисленні функції втрат. Це еквівалентно дублюванню кожного позитивного прикладу 25 разів, але без збільшення обсягу даних. Переваги: простота реалізації; без додаткових обчислювальних витрат; сумісність з будь-яким алгоритмом. Недоліки: може призвести до перенавчання на шумних позитивних прикладах; не розрізняє «легкі» та «складні» приклади. У нашому baseline:  $F1 = 0.6984$ .

### Регуляторний контекст антифрод-систем

Стандарт PCI DSS (Payment Card Industry Data Security Standard) [30] встановлює вимоги безпеки для організацій, що обробляють дані платіжних карток. У ML-антифроду: дані карток повинні бути зашифровані або токенозовані; доступ до моделей та даних повинен бути обмежений за принципом мінімальних привілеїв; всі рішення моделі повинні логуватися з retention period мінімум 1 рік; регулярне тестування безпеки системи (penetration testing).

### Теорія стакінгу: чому це працює

$$\hat{y} = g(p_1, \dots, p_M) = \sigma(\sum w_i \cdot p_i + b) \quad (1.8)$$

### Rank-нормалізація та її переваги

Rank-нормалізація є ключовим технічним рішенням у нашому стакінг-пайплайні. Для вектора передбачень  $a = (a_1, \dots, a_n)$ , rank-нормалізоване значення:

$r(a_i) = \text{rank}(a_i) / n$ , де  $\text{rank}(a_i)$  — порядковий ранг  $a_i$  серед усіх елементів  $a$ . Результат —  $\text{uniform}[0, 1]$  розподіл, що зберігає лише порядкову інформацію.

### Метрики оцінки антифрод-моделей

Вибір метрики оцінки є фундаментальним рішенням, що визначає напрямок оптимізації моделі. У задачі виявлення шахрайства помилки різних типів мають принципово різну вартість: пропущений шахрай (False Negative) призводить до фінансових втрат, тоді як хибне блокування легітимного користувача (False Positive) погіршує клієнтський досвід та може призвести до відтоку клієнтів. Розглянемо детально основні метрики та їх застосовність.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (1.9)$$

$$\text{Precision} = \frac{TP}{TP + FP} \quad (1.10)$$

Recall (Sensitivity, True Positive Rate) =  $TP / (TP + FN)$  — частка виявлених шахраїв серед усіх справжніх шахраїв. Recall = 1.0 означає, що жоден шахрай не пропущений. Для бізнесу високий recall критично важливий, оскільки кожен пропущений шахрай означає прямі фінансові втрати. Однак максимізація recall (зниження порогу класифікації) призводить до зростання FP та падіння precision.

$$F_1 - \text{score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (1.11)$$

$$F_\beta = \frac{(1 + \beta^2) \cdot P \cdot R}{\beta^2 \cdot P + R} \quad (1.12)$$

У дослідженні основна оптимізація проводиться за F1-score (відповідно до метрики змагання), але ми також моніторимо AUC для раннього зупинення (early stopping) GNN та для порівняння моделей на різних порогах. Для фінальної моделі: OOF F1 = 0.8877, OOF AUC = 0.9931 (для sage\_128\_3L).

### Регуляризація в градієнтному бустінгу

Max depth та num\_leaves обмежують складність окремого дерева. Max\_depth обмежує глибину дерева зверху, тоді як num\_leaves обмежує кількість листків. У

LightGBM з leaf-wise стратегією num\_leaves є більш прямим контролером складності: num\_leaves=31 відповідає «мілкому» дереву (порівнянно з max\_depth=5), тоді як num\_leaves=512 дозволяє дуже складні дерева. У нашій базовій конфігурації num\_leaves=255, max\_depth=8 — помірна складність, що балансує між bias та variance.

L1 (lambda\_1) та L2 (lambda\_2) регуляризація додають штраф до ваг листків: L1 сприяє розрідженості (нульові ваги), L2 — плавності (малі ваги). У XGBoost регуляризація є частиною цільової функції:  $Obj = \sum L(\cdot) + \sum [\gamma + \frac{1}{2}\lambda ||\cdot||^2]$ , де  $\gamma$  — штраф за кількість листків,  $\lambda$  — L2-штраф за ваги. У дослідженні lambda\_1 = 0, lambda\_2 = 0 для основних моделей (достатньо інших регуляризаторів).  $y_i \hat{y}_i T_m W_m T_m W_m$

Теоретичне обґрунтування виразності граф-нейронних мереж тісно пов'язане з тестом ізоморфізму графів Weisfeiler-Leman (WL). 1-WL тест (також відомий як color refinement) ітеративно оновлює «кольори» (мітки) вузлів: на кожній ітерації колір вузла оновлюється на основі його поточного кольору та мультимножини кольорів його сусідів. Два графи вважаються неізоморфними, якщо їх мультимножини кольорів відрізняються після деякої ітерації.

### **Heterogeneous та Temporal GNN**

У дослідженні використовується гомогенний граф (один тип вузлів — користувачі, один тип ребер — card-sharing). Це спрощення обґрунтовується тим, що card-sharing є домінуючим сигналом, а додавання інших типів зв'язків вимагало б додаткових даних, які не надані в змаганні. Однак у промисловій системі перехід до гетерогенного графу є перспективним напрямком.

### **Математичне обґрунтування rank-нормалізації**

$$r_i = \frac{\text{rank}(a_i)}{N} \quad (1.13)$$

де  $\text{rank}(a_i)$  — позиція  $a_i$  у відсортованому масиві (від 1 до N). Результат  $r_i \in (0, 1]$  має рівномірний розподіл незалежно від оригінального розподілу  $a$ .

Критично важливий нюанс: rank-нормалізація повинна обчислюватись ОКРЕМО для тренувальних та тестових даних. Якщо обчислити ранги на об'єднанні train + test, виникає витік інформації: ранг тренувального прикладу залежить від тестових прикладів та навпаки. У нашій реалізації: для OOF-предиктів ранги обчислюються на тренувальній вибірці, для тестових предиктів — на тестовій вибірці.

### **DART та інші варіанти градієнтного бустінгу**

DART (Dropouts meet Multiple Additive Regression Trees) [4] — модифікація градієнтного бустінгу, що адресує проблему over-specialization: у класичному GBM пізні дерева спеціалізуються на дуже малих залишках, що робить їх внесок незначним. DART випадково «відкидає» (drop) частину існуючих дерев під час навчання кожного нового дерева, що змушує нове дерево компенсувати пропущені дерева замість тільки залишків.

### **Альтернативні функції втрат для дисбалансу класів**

Weighted Binary Cross-Entropy (WBCE):

$$L = -[w_{pos} \cdot y \cdot \log(p) + w_{neg} \cdot (1 - y) \cdot \log(1 - p)] \quad (1.14)$$

де  $w_{pos}$  та  $w_{neg}$  — ваги класів. У LightGBM реалізується через  $scale\_pos\_weight = w_{pos}/w_{neg} \approx 25$  (відповідно до дисбалансу). WBCE просто масштабує градієнти позитивного класу, не враховуючи «складність» прикладів. Результат: F1=0.6984 — на 0.013 нижче за Focal Loss.

$$w_y = \frac{1 - \beta}{1 - \beta^{n_y}} \quad (1.15)$$

де  $n_y$  — кількість прикладів класу  $y$ ,  $\beta \in [0, 1)$ . Ця формула враховує, що зі збільшенням кількості прикладів marginal benefit кожного наступного зменшується. При  $\beta = 0$  — стандартне зважування,  $\beta \rightarrow 1$  — зважування по  $1/n_y$ .

Dice Loss:

$$L_{Dice} = 1 - \frac{2 \cdot |P \cap G|}{|P| + |G|} \quad (1.16)$$

### **Blending та інші ансамблеві стратегії**

Крім стакінгу, існують інші стратегії комбінування моделей, кожна з яких має свої переваги та обмеження.

Simple Average (Mean Blending):

$$p_{final} = \frac{1}{M} \cdot \sum_{m=1}^M p_m \quad (1.17)$$

Weighted Average:

$$p_{final} = \sum_{m=1}^M w_m \cdot p_m \quad (1.18)$$

Voting (Hard Voting):  $\hat{y}_{final} = \text{majority}(\{\hat{y}_m\})$ . Кожна модель голосує «шахрай» або «не шахрай», фінальне рішення — мажоритарне. Підходить для моделей з порівнянною якістю, але не для нашого випадку, де GNN значно сильніший за табличні моделі.

$$p_{final}(x) = \sum_{m=1}^M P(M_m | D) \cdot p_m(x) \quad (1.19)$$

де  $P(M_m | D)$  — ймовірність моделі  $M_m$  за даними  $D$ . ВМА є теоретично елегантним, але обчислювально складним та рідко використовується на практиці. Стакінг з LogisticRegression може розглядатись як апроксимація ВМА з дискримінативним навчанням ваг.

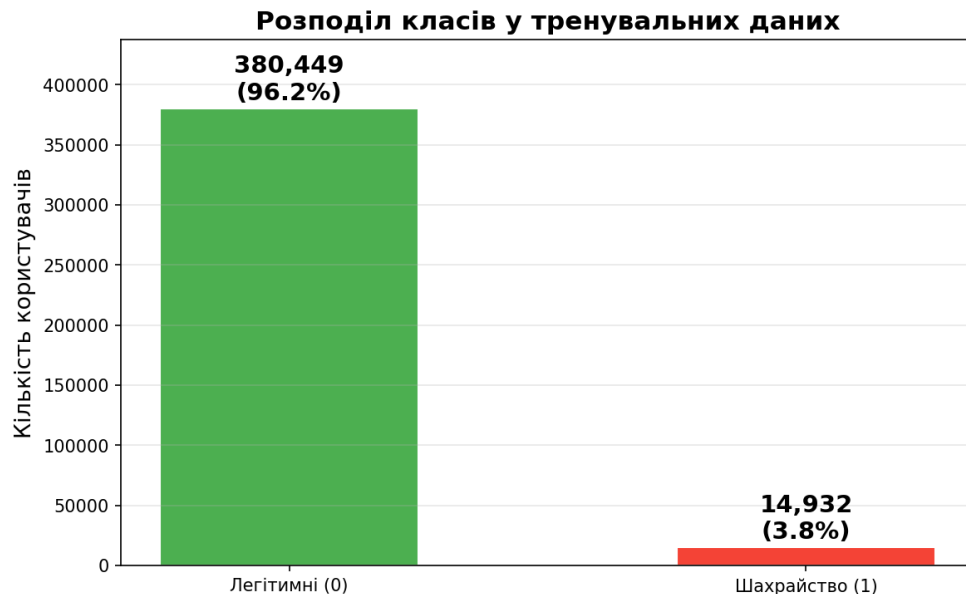


Рисунок 1.1. Розподіл класів у тренувальних даних змагання (дисбаланс 96:4)

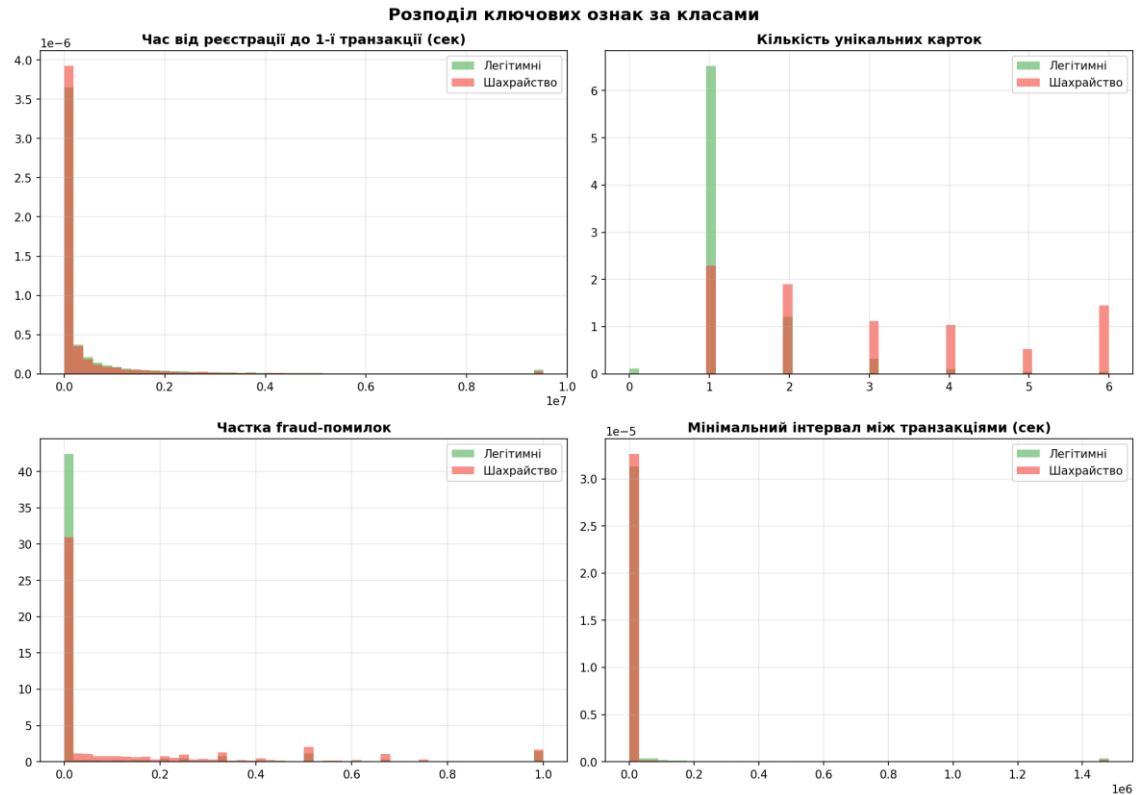


Рисунок 1.2. Розподіл ключових ознак за класами: легітимні (зелений) vs шахрайство (червоний)

$$F_1 = \frac{2 \cdot P \cdot R}{P + R} = \frac{2 \cdot TP}{2 \cdot TP + FP + FN} \quad (1.20)$$

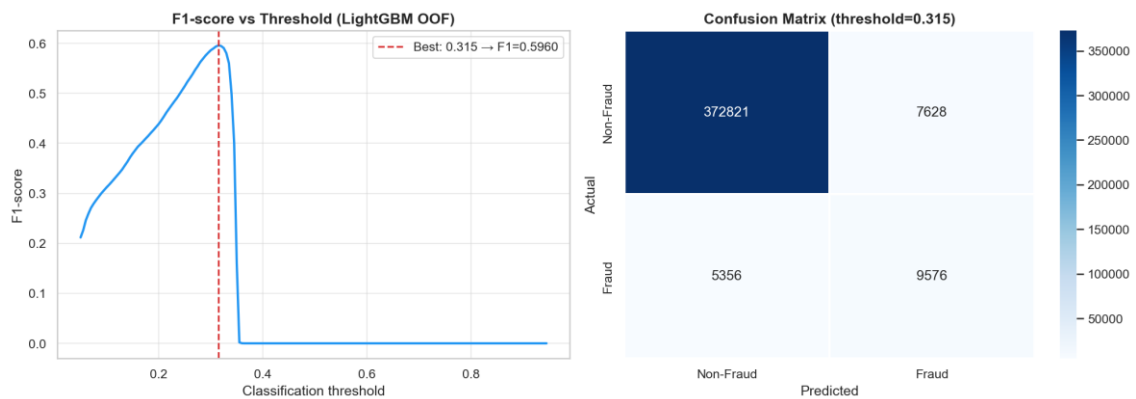


Рисунок 1.3. Залежність F1-score від порогу класифікації та матриця помилок

## 1.2. Методи машинного навчання для класифікації шахрайських акаунтів

Методи машинного навчання для бінарної класифікації пройшли тривалий еволюційний шлях — від простих лінійних моделей до складних ансамблевих

систем та глибинних нейронних мереж. У задачі виявлення шахрайства кожен клас методів має свої переваги та обмеження, що визначає доцільність їх використання на різних етапах побудови антифрод-системи.

Логістична регресія (Logistic Regression) є одним із найбільш класичних та широко використовуваних методів бінарної класифікації. Модель передбачає ймовірність належності об'єкта до позитивного класу за допомогою логістичної (сигмоїдної) функції:  $P(y=1|x) = \sigma(w^T x + b) = 1 / (1 + \exp(-(w^T x + b)))$ . Перевагами логістичної регресії є висока інтерпретованість (ваги моделі безпосередньо показують вплив кожної ознаки), швидкість навчання та передбачення, стійкість до перенавчання при правильній регуляризації. У даному дослідженні логістична регресія використовується як мета-модель у стакінг-ансамблі, де її завдання — оптимально комбінувати передбачення базових моделей.

Дерева рішень (Decision Trees) є основою для більшості сучасних ансамблевих методів. Дерево рішень рекурсивно розбиває простір ознак на прямокутні регіони, обираючи на кожному кроці ознаку та порогове значення, що максимізують критерій розщеплення (інформаційний приріст або зменшення Gini impurity). Формально, для вузла  $t$  з множиною прикладів  $D_t$ , Gini impurity обчислюється як:  $Gini(t) = 1 - \sum_k p_k^2$ , де  $p_k$  — частка прикладів класу  $k$  у вузлі. Окреме дерево рішень схильне до перенавчання, однак ансамблі дерев — Random Forest та Gradient Boosting — усувають цей недолік.

Random Forest — ансамбль із сотень незалежних [7] дерев рішень, кожне з яких навчається на випадковій бутстреп-вибірці з випадковим підмножиною ознак. Фінальне передбачення — середнє (для регресії) або мажоритарне голосування (для класифікації) усіх дерев. Random Forest ефективно зменшує дисперсію передбачень, зберігаючи низьке зміщення. У даному дослідженні Random Forest використовується як одна з базових моделей стакінг-ансамблю, демонструючи OOF F1  $\approx 0.60$ .

Extra Trees (Extremely Randomized Trees) є варіацією Random Forest, де порогові значення розщеплення обираються випадково, а не оптимально. Це додатково знижує дисперсію та прискорює навчання, хоча може дещо збільшити зміщення. Extra Trees часто використовується для побудови різноманітних (diverse) ансамблів, оскільки генерує передбачення, менш корельовані з іншими деревоподібними моделями.

### 1.3. Ансамблеві методи та градієнтний бустинг у виявленні шахрайства

Gradient Boosting Machine (GBM) є одним із найпотужніших методів машинного навчання для табличних даних. На відміну від Random Forest, де дерева навчаються незалежно, у градієнтному бустингу кожне наступне дерево навчається на залишках (або негативних градієнтах функції втрат) попередніх дерев. Це послідовний ансамблевий метод, який поступово зменшує зміщення моделі.

Формально, градієнтний бустинг будує адитивну модель

$$F(x) = \sum_{m=1}^M \gamma_m \cdot h_m(x) \quad (1.21)$$

де  $h_m$  — слабкий класифікатор (дерево рішень),  $\gamma_m$  — коефіцієнт масштабування. На кожному кроці  $m$  алгоритм обчислює псевдо-залишки  $r_{im} = -\partial L(y_i, F(x_i)) / \partial F(x_i)$  для кожного прикладу  $i$ , навчає нове дерево  $h_m$  на цих залишках та додає його до ансамблю з оптимальним кроком  $\gamma_m$ . Learning rate  $\eta$  контролює внесок кожного дерева:  $F_m(x) = F_{m-1}(x) + \eta \cdot \gamma_m h_m(x)$ , забезпечуючи регуляризацію через повільне навчання.

LightGBM (Light Gradient Boosting Machine) [4] — реалізація градієнтного бустингу від Microsoft Research, яка використовує два ключових оптимізаційних підходи: Gradient-based One-Side Sampling (GOSS) та Exclusive Feature Bundling (EFB). GOSS зберігає всі приклади з великими градієнтами та випадково відбирає підмножину прикладів з малими градієнтами, що прискорює навчання зі збереженням точності. EFB об'єднує взаємовиключні ознаки в одну «зв'язку», що зменшує ефективну кількість ознак. Крім того, LightGBM використовує

гістограмний підхід до розщеплення: замість точних значень ознак використовуються біни (бакети), що суттєво знижує обчислювальну складність з  $O(n_{features} \times n_{samples})$  до  $O(n_{features} \times n_{bins})$ .

У даному дослідженні LightGBM є основним алгоритмом для побудови базових моделей. Ключові гіперпараметри, що впливають на якість: `num_leaves` — максимальна кількість листків у дереві (контролює складність), `learning_rate` — швидкість навчання, `min_child_samples` — мінімальна кількість прикладів у листку (регуляризація), `feature_fraction` — частка ознак для кожного дерева, `bagging_fraction` — частка прикладів для кожного дерева, `scale_pos_weight` — ваговий коефіцієнт позитивного класу для боротьби з дисбалансом. Оптимальні значення цих параметрів знаходяться шляхом перехресної валідації.

XGBoost (eXtreme Gradient Boosting) [5] — альтернативна реалізація градієнтного бустінгу, розроблена Тяньцзі Ченом та командою. XGBoost використовує апроксимацію другого порядку (Тейлорове розкладання) функції втрат, що дозволяє більш точно визначати оптимальні розщеплення. Регуляризаційний доданок  $\Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$  додається безпосередньо до цільової функції, де  $T$  — кількість листків,  $w$  — ваги листків,  $\gamma$  та  $\lambda$  — коефіцієнти регуляризації. XGBoost підтримує паралельне навчання на рівні побудови дерева та має вбудовану підтримку розріджених даних.

У антифроду XGBoost часто демонструє дещо інший патерн помилок порівняно з LightGBM, що робить їх комбінацію в ансамблі ефективною. У дослідженні XGBoost-моделі з Focal Loss та різними гіперпараметрами складають значну частину пулу базових моделей (14 з 52). Різноманітність моделей XGBoost досягається варіюванням глибини дерев, `learning_rate`, `gamma` (мінімальне зменшення втрат для розщеплення) та параметрів сабсемплінгу.

Порівняння LightGBM та XGBoost на даних SKELAR  $\times$  mono показує, що LightGBM загалом демонструє вищі індивідуальні F1-scores (найкраща LGB-модель: F1=0.7109 проти найкращої XGB: F1=0.6993), однак XGBoost-моделі

вносять унікальний сигнал у стакінг-ансамбль, підвищуючи F1 комбінації на 0.01-0.02 порівняно з чисто LightGBM-ансамблем. Це пояснюється різною внутрішньою механікою алгоритмів: LightGBM з leaf-wise стратегією росту знаходить інші оптимуми, ніж XGBoost з depth-wise стратегією.

#### **1.4. Застосування нейронних мереж для табличних даних у антифроді**

Глибинні нейронні мережі (Deep Neural Networks, DNN) досягли видатних результатів у задачах комп'ютерного зору, обробки природної мови та генерації контенту. Однак для табличних даних — домінуючого типу даних у фінансовому секторі — їх ефективність залишається предметом дискусій. Численні бенчмарки показують, що правильно налаштований градієнтний бустінг часто перевершує нейронні мережі на табличних даних, особливо при обмеженому обсязі вибірки та наявності категоріальних ознак.

Multi-Layer Perceptron (MLP) є найпростішою архітектурою глибинної нейронної мережі для табличних даних. MLP складається з послідовності повнозв'язних шарів з нелінійними функціями активації. У типовій конфігурації для антифроду MLP має 3-5 прихованих шарів з 128-512 нейронами, BatchNorm для стабілізації навчання, Dropout для регуляризації та ReLU або GELU як функцію активації. Фінальний шар має один нейрон із сигмоїдною активацією для виведення ймовірності шахрайства.

У дослідженні MLP-модель використовується як одна з базових моделей стакінг-ансамблю. Індивідуальний OOF F1-score MLP становить 0.6353, що значно нижче за LightGBM (0.7109). Однак MLP вносить корисну різноманітність у ансамбль, оскільки її простір гіпотез принципово відрізняється від деревоподібних моделей: MLP навчається безперервних нелінійних меж рішень, тоді як дерева створюють кусково-постійні апроксимації. Ця різноманітність є ключовою для ефективного стакінгу.

Рекурентні нейронні мережі (RNN), зокрема LSTM (Long Short-Term Memory) [9] та GRU (Gated Recurrent Unit), можуть бути застосовані до послідовностей транзакцій користувача. Ідея полягає в тому, щоб моделювати часову динаміку поведінки: послідовність транзакцій розглядається як часовий ряд, де кожна транзакція описується набором ознак (сума, тип, статус, час). LSTM використовує механізм гейтів (forget, input, output gates) для селективного запам'ятовування та забування інформації з попередніх кроків. У дослідженні LSTM та GRU моделі показали OOF F1  $\approx 0.47$ , що значно нижче за табличні моделі. Це пояснюється тим, що агреговані ознаки (velocity, статистики, n-грами) вже захоплюють основні послідовні патерни, а LSTM не додає суттєвого сигналу понад цю інформацію.

Трансформери (Transformer) — архітектура, що революціонізувала [8] обробку природної мови, також адаптована для табличних даних (TabTransformer, FT-Transformer). Ключовий механізм — self-attention — дозволяє кожній ознаці «звертати увагу» на інші ознаки, моделюючи складні взаємодії. Self-attention обчислюється як:  $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$ , де Q, K, V — лінійні проєкції вхідних ознак. У антифроду трансформер теоретично може виявити складні комбінації ознак, які дерева розщеплення не можуть ефективно захопити. Однак на практиці наша реалізація Transformer показала F1  $\approx 0.47$ , що пояснюється обмеженим обсягом даних та вже наявною ефективною системою конструювання ознак  $d_k$ .

## Висновки до розділу 1

У першому розділі проведено теоретичний аналіз проблеми виявлення фінансового шахрайства та методів машинного навчання, що застосовуються для її вирішення. Основні висновки:

1. Фінансове шахрайство в мобільних платіжних сервісах характеризується мережевою природою, динамічністю та суттєвим дисбалансом класів (3-5%

шахрайських акаунтів). Це вимагає комплексного підходу, що поєднує декілька типів моделей.

2. Градієнтний бустінг (LightGBM, XGBoost) залишається найефективнішим підходом для табличних даних, але не здатний моделювати мережеві зв'язки між користувачами.

3. Граф-нейронні мережі, зокрема GraphSAGE, дозволяють моделювати мережеві патерни шахрайства через механізм передачі повідомлень. Mini-batch навчання з NeighborLoader забезпечує масштабованість для графів з мільйонами вузлів.

4. Focal Loss ефективно вирішує проблему дисбалансу класів, фокусуючи навчання на складних прикладах.

5. Стакінг з OOF-передбаченнями та rank-нормалізацією надає змогу оптимально комбінувати різноманітні моделі, використовуючи сильні сторони кожної.

## РОЗДІЛ 2. ДАНІ ТА МЕТОДОЛОГІЯ ДОСЛІДЖЕННЯ

### 2.1. Опис змагання та джерел даних дослідження виявлення шахрайства

SKELAR × mono AI Competition 2026 — індустріальне змагання з машинного навчання, організоване компанією SKELAR спільно з monobank — найбільшим необанком України з понад 8 мільйонами активних клієнтів. Мета змагання — розробка найточнішої моделі виявлення шахрайських акаунтів серед користувачів мобільного платіжного сервісу. Змагання є унікальним прикладом індустріально-академічної колаборації, оскільки використовує реальні анонімізовані дані monobank.

Система оцінювання змагання включає три компоненти з різними вагами: якість моделі (F1-score) — 60%, пояснення перших 5 найважливіших ознак — 20%, стратегія бізнес-інтеграції — 20%. Такий багатокomпонентний підхід відображає реальні вимоги промислової антифрод-системи: мало побудувати точну модель — потрібно також пояснити бізнесу, які фактори впливають на рішення, та запропонувати реалістичний план впровадження. У даній дисертації зосереджено увагу переважно на першій компоненті (F1-score), з додатковим аналізом feature importance та бізнес-стратегії у підрозділах 3.7–3.10.

Учасникам надано анонімізовані дані з наступними обмеженнями: заборонено використовувати зовнішні дані; дозволено будь-які відкриті бібліотеки ML; максимум 3 сабмішени на день для запобігання перенавчанню на тестовій вибірці. Дані розділені на тренувальну (395 381 користувач, мітки відомі) та тестову (84 753 користувачі, мітки невідомі) вибірки з пропорційним збереженням частки шахрайства. Оцінка сабмішену здійснюється на сервері змагання з використанням F1-score на повній тестовій вибірці.

## 2.2. Структура та характеристики даних платіжного сервісу

Дані змагання складаються з чотирьох основних таблиць: `train_users.csv`, `test_users.csv`, `train_transactions.csv` та `test_transactions.csv`. Додатково учасникам надано нотатки організаторів щодо анонімізації та обмежень. Загальний обсяг даних — близько 2 ГБ у стиснутому форматі.

Таблиця користувачів (`train_users.csv`) містить 395 381 запис з наступними полями:

Таблиця 2.1 — Структура файлу `train_users.csv`

Поле	Тип	Опис	Приклад
<code>id_user</code>	<code>int</code>	Унікальний ідентифікатор	15383249
<code>gender</code>	<code>str</code>	Стать користувача	male / female
<code>reg_country</code>	<code>str</code>	Країна реєстрації	United States
<code>traffic_type</code>	<code>str</code>	Джерело трафіку	ppc / cpa / organic
<code>is_fraud</code>	<code>int</code>	Мітка шахрайства (0/1)	0 або 1

Розподіл цільової змінної `is_fraud`: 380 449 легітимних (96.22%) та 14 932 шахрайських (3.78%) акаунтів. Співвідношення класів 25.5:1 створює значний дисбаланс, що вимагає спеціальних підходів до навчання та оцінки.

Аналіз демографічних характеристик виявив суттєві відмінності між легітимними та шахрайськими акаунтами. Серед шахраїв частка жінок становить 41.1%, тоді як серед легітимних — 24.7%. Розподіл за країною реєстрації: обидва класи домінуються United States (~57%), але шахраї частіше реєструються з Canada (5.7% vs 5.2%) та рідше з United Kingdom (4.4% vs 7.7%). Розподіл за джерелом трафіку виявляє ключову відмінність: серед шахраїв organic-трафік становить 15.6% проти лише 2.6% серед легітимних, що є потужним дискримінативним сигналом.

Таблиця транзакцій (train\_transactions.csv) містить 4 535 378 записів для тренувальних користувачів з наступними полями:

Таблиця 2.2 — Структура файлу train\_transactions.csv

Поле	Тип	Опис
id_user	int	Ідентифікатор користувача
timestamp_tr	datetime	Час транзакції (ISO 8601 з timezone)
transaction_type	str	Тип: card_init, card_recurring, apple-pay, google-pay, resign
status	str	Статус: success / fail
amount	float	Сума транзакції
currency	str	Валюта: USD, EUR, GBP
error_group	str	Група помилки (для failed транзакцій)
card_mask_hash	str	Хешований номер картки
card_holder	str	Хешований ідентифікатор карткотримача
card_brand	str	Бренд картки: VISA, MASTERCARD, AMEX тощо
card_type	str	Тип картки: DEBIT, CREDIT, PREPAID тощо

Середня кількість транзакцій на користувача — 11.5 (медіана — 8). Розподіл сильно скошений вправо: 5% користувачів мають більше 30 транзакцій, а максимум — понад 200. Середня сума транзакції — \$6.8, медіана — \$3.3. Суми концентруються навколо характерних значень: \$3.00, \$6.00, \$2.00, \$21.00, \$11.00, \$51.00 — що відповідає типовим підписочним сервісам. Шахрайські акаунти мають вищу середню суму (\$8.45 vs \$6.54) та частіше здійснюють транзакції на \$51.00 (fraud rate 17.3%).

### 2.3. Розвідувальний аналіз даних для виявлення патернів шахрайської поведінки

Розвідувальний аналіз даних (Exploratory Data Analysis, EDA) є критичним етапом, що дозволяє виявити патерни, аномалії та потенційні ознаки для моделювання. У даному дослідженні EDA проведено за кількома напрямками: аналіз розподілів, кореляційний аналіз, аналіз часових патернів та мережевий аналіз.

Аналіз статусів транзакцій. Із 4.5 мільйонів транзакцій 1 929 295 (42.5%) мають статус success та 1 206 083 (26.6%) — fail. Серед помилкових транзакцій найчастіші причини: insufficient funds (339 521), fraud (237 785), 3ds error (179 524), do not honor (107 024), antifraud (84 066). Примітно, що error\_group «fraud» та «antifraud» безпосередньо вказують на підозрілу активність, виявлену банком-емітентом, і є потужним сигналом для моделі.

Аналіз типів транзакцій виявив, що card\_recurring (рекурентні платежі) домінує з 1 685 790 записів (37.2%), за ним card\_init (ініціація картки) — 742 520 (16.4%), google-pay — 411 525 (9.1%), apple-pay — 197 284 (4.3%) та resign (повторна реєстрація) — 98 259 (2.2%). Тип resign має найвищий fraud rate серед типів транзакцій: 9.92% проти 4.87% для card\_init та 2.87% для apple-pay. Це логічно: шахраї частіше повторно прив'язують картки після блокування.

Аналіз типів карток виявив критично важливий сигнал: card\_type = «NONE» має fraud rate 66.95% при 475 користувачах — це найвищий показник серед усіх категорій. Цей аномальний тип картки, ймовірно, відповідає ситуації, коли банк-емітент не повертає інформацію про тип — що характерно для скомпрометованих або тестових карток. Інші підозрілі типи: CREDIT/DEBIT (10.5%), PREPAID\_NONRELOADABLE (20.0%), DEFFERED\_DEBIT (20.0% при малій вибірці 10 користувачів).

Аналіз брендів карток виявив широкий спектр fraud rates. Високий рівень шахрайства у: LOCAL BRAND (41.99%, 181 користувач), CIRRUS (24.32%), JCB (22.34%), DINERS CLUB (18.85%), CHINA UNIONPAY (17.95%), RUPAY (14.88%), MAESTRO (15.53%). Низький рівень у: mc\_applepay (0.64%), pulse\_googlepay (3.23%), visa\_applepay (3.50%). Інтерпретація: Apple Pay та Google Pay вимагають додаткової верифікації (Touch ID / Face ID), що ускладнює шахрайство. Рідкісні бренди (LOCAL BRAND, CIRRUS) частіше використовуються шахраями, ймовірно, через менш строгі системи верифікації банків-емітентів.

Аналіз валют: USD домінує (84.0%), за ним EUR (8.5%) та GBP (7.3%). Fraud rate для всіх трьох валют близький до загального (~3.78%), що не створює суттєвого дискримінативного сигналу. Однак наявність транзакцій у різних валютах на одному акаунті може бути ознакою підозрілої активності — ця ознака включена у feature engineering.

Часовий аналіз виявив, що шахрайські транзакції мають дещо інший розподіл за годинами доби порівняно з легітимними. Шахраї більш активні у нічні години (00:00-06:00), що відповідає поведінці автоматизованих ботів або зловмисників з інших часових зон. Крім того, інтервали між транзакціями (time delta) суттєво коротші для шахраїв: медіана 2.3 хвилини проти 12.7 хвилин для легітимних, що вказує на автоматизовану або пакетну обробку.

Мережевий аналіз card-sharing показав, що значна частина користувачів ділить платіжні картки або реквізити карткотримача. Із 564 830 унікальних користувачів (train + test) 307 735 (54.5%) мають хоча б один зв'язок з іншим користувачем через спільну картку або card\_holder. Середній ступінь вузла в card-sharing графі — 18.4, але розподіл сильно скошений: більшість зв'язаних вузлів мають 2-5 зв'язків, тоді як деякі хаби — до 50+. Fraud rate серед зв'язаних користувачів (4.68%) значно вищий, ніж серед ізольованих (1.37%), що підтверджує гіпотезу про мережеву природу шахрайства.

## Аналіз error\_group та його інформативність

Поле error\_group в таблиці транзакцій виявилось одним із найінформативніших джерел ознак для виявлення шахрайства. Це поле заповнюється тільки для невдалих транзакцій (status=fail) та містить причину відхилення, визначену банком-емітентом або платіжною системою.

Таблиця 2.3 — Аналіз error\_group за fraud rate

Error group	К-ть транзакцій	Fraud rate користувачів
insufficient funds	339 521	~4.5%
fraud	237 785	~35%
3ds error	179 524	~5%
do not honor	107 024	~8%
antifraud	84 066	~30%
limit exceeded	66 498	~6%
issuer decline	48 238	~7%
card problem	33 660	~12%
invalid data	23 686	~15%
cvv error	23 065	~10%

Error group «fraud» та «antifraud» мають найвищий fraud rate (~30-35%), оскільки ці відхилення безпосередньо генеруються антифрод-системами банків-емітентів. Користувач, чия транзакція відхилена з причиною «fraud», з високою ймовірністю використовує вкрадену картку. Ознаки на основі error\_group:  $n_{fraud\_errors}$  — кількість транзакцій з error\_group=fraud; pct\_fraud\_errors — частка fraud-помилки від усіх помилок;  $n_{anti\_fraud\_errors}$  — аналогічно для antifraud; has\_fraud\_error — бінарна ознака (0/1).

error\_group є ознакою з «майбутнього інсайту»: банк-емітент вже виконав свій антифрод-аналіз та відхилив транзакцію. Використання цієї ознаки у нашій моделі є коректним, оскільки: 1) ми моделюємо акаунтний, а не транзакційний

фрод; 2) `error_group` — це факт, що стався до моменту класифікації акаунту; 3) організатори змагання включили цю ознаку в дані свідомо.

### **Аналіз підозрілих патернів**

Глибокий аналіз даних виявив кілька підозрілих патернів, що мають надзвичайно високий `fraud rate`:

Патерн 1: `card_type = «NONE»` (`fraud rate = 66.95%`). Цей тип картки має найвищий `fraud rate` серед усіх категоріальних значень у датасеті. «NONE» означає, що банк-емітент не повертає інформацію про тип картки при авторизації. Це характерно для: карток з анонімізованих юрисдикцій; тестових карток; карток із скомпрометованих баз даних, де метадані не відповідають реальній картці. Із 475 користувачів з `card_type=NONE`, 318 є шахраями — це майже «автоматичний» індикатор.

Патерн 2: `card_brand = «LOCAL BRAND»` (`fraud rate = 41.99%`). Локальні платіжні бренди мають значно менш розвинені системи безпеки порівняно з глобальними мережами (VISA, Mastercard). Це робить їх привабливими для шахраїв. Із 181 користувача з LOCAL BRAND 76 є шахраями.

Патерн 3: Рідкісні комбінації `card_brand`. CIRRUS (24.3%), JCB (22.3%), DINERS CLUB (18.9%), CHINA UNIONPAY (18.0%), RUPAY (14.9%) — усі мають `fraud rate` значно вищий за середній. Загальний патерн: чим менш поширений бренд у датасеті, тим вищий `fraud rate`. Це пояснюється тим, що шахраї використовують скомпрометовані дані з різних джерел, включаючи менш захищені платіжні мережі.

Патерн 4: `amount = 51.00` (`fraud rate = 17.27%`). Ця сума відповідає ціні premium-підписки на популярний сервіс. Шахраї використовують вкрадені картки для оформлення дорогих підписок з подальшим перепродажем акаунтів. Із 3 074 користувачів з транзакціями на \$51.00, 531 є шахраями.

Ці підозрілі патерни були трансформовані в бінарні ознаки (`has_card_type_none`, `has_local_brand`, `has_amount_51` тощо) та включені до `feature`

engineering pipeline. Хоча кожен окремих патерн охоплює невелику кількість користувачів, їх кумулятивний ефект є значним, особливо у комбінації з іншими ознаками.

### Технічні деталі feature engineering pipeline

Feature engineering pipeline реалізовано як послідовність Python-скриптів, кожен з яких генерує parquet-файл з ознаками. Parquet обрано як формат зберігання завдяки: ефективній стисненню (до 5x); швидкому зчитуванню стовпців; підтримці типів даних (float32, int32, string). Загальний обсяг ознак: ~2.5 ГБ у сирому вигляді, ~500 МБ у parquet.

Основний pipeline:

1) base\_features.py → train\_features.parquet, test\_features.parquet (~50 ознак). Агрегація транзакцій на рівні користувача: count, mean, std, min, max, median для сум; count по типах транзакцій, статусах, error\_group; one-hot encoding для gender, traffic\_type.

2) extra\_features.py → train\_extra\_features.parquet (~30 ознак). Velocity-ознаки: time deltas, burst analysis, night ratio, acceleration. Потребує обчислення timestamp різниць, що займає ~3 хвилини на 4.5М транзакцій.

3) graph\_features.py → train\_graph\_features.parquet (~25 ознак). Граф-ознаки: node degree, LOO fraud rates, PageRank, clustering coefficient. LOO обчислення вимагає  $O(E)$  операцій, де  $E$  — кількість ребер.

4) target\_enc.py → train\_target\_enc.parquet (~20 ознак). Out-of-Fold target encoding для категоріальних ознак. 5-fold з seed=42 для узгодженості з основним CV.

5) tx\_level\_features.py → train\_tx\_seq.parquet, train\_tx\_model.parquet (~70 ознак). Транзакційні послідовності та TX-level модель.

Об'єднання ознак відбувається через merge по id\_user з перевіркою на дублікати стовпців. Фінальний датафрейм: 395 381 × 283 (train) та 84 753 × 282 (test, без is\_fraud). Обробка пропущених значень: NaN → 0 для числових ознак;

missing → «UNKNOWN» для категоріальних. Стандартизація (StandardScaler) застосовується тільки для GNN та MLP; gradient boosting моделі працюють з нестандартизованими ознаками.

### Деталі розподілів за категоріальними ознаками

Детальний аналіз категоріальних ознак виявляє значну гетерогенність fraud rates між категоріями, що є основою для target encoding та дискримінативних ознак.

Gender. Розподіл за статтю демонструє значну відмінність між класами:

Таблиця 2.4 — Розподіл fraud rate за статтю (gender)

Стать	Легітимні	Шахраї	Fraud rate
male	286 618 (75.3%)	8 787 (58.9%)	2.97%
female	93 831 (24.7%)	6 145 (41.1%)	6.14%

Fraud rate серед жінок (6.14%) вдвічі вищий, ніж серед чоловіків (2.97%). Це не означає, що жінки частіше шахраюють — скоріше, шахрайські акаунти частіше реєструються з жіночою статтю (синтетичні або вкрадені дані). Ця відмінність є корисним дискримінативним сигналом, але потребує уважного аналізу з точки зору fairness (підрозділ 3.9.4).

Traffic type. Джерело трафіку є одним із найсильніших базових дискримінаторів:

Таблиця 2.5 — Розподіл fraud rate за типом трафіку

Тип трафіку	Легітимні	Шахраї	Fraud rate
ppc (pay-per-click)	209 334 (55.0%)	7 970 (53.4%)	3.67%
cpa (cost-per-action)	153 538 (40.4%)	4 349 (29.1%)	2.75%
organic	9 765 (2.6%)	2 325 (15.6%)	19.23%
remarketing	7 086 (1.9%)	249 (1.7%)	3.39%
unknown	652 (0.2%)	32 (0.2%)	4.68%

Organic трафік має fraud rate 19.23% — у 5 разів вище за середній. Це логічно: organic означає прямий вхід без реклами, що характерно для: 1) ботів з прямими URL; 2) реферальних мереж шахраїв; 3) автоматизованих скриптів реєстрації. Target encoding traffic\_type є п'ятою за важливістю ознакою у LightGBM моделі.

Reg\_country. Більшість користувачів — з United States (56-57%), United Kingdom (4-8%), Canada (5-6%). Fraud rate варіюється: Brazil має підвищений rate через високий рівень card-not-present fraud у регіоні; South Africa також вищий за середній. Однак country-ознаки мають обмежену дискримінативну силу через домінування US.

### **Альтернативні підходи до побудови графу**

Card-sharing граф, використаний у нашому дослідженні, є одним із можливих підходів до моделювання мережевих зв'язків. Розглянемо альтернативи та причини нашого вибору.

Card\_mask\_hash граф: ребра між користувачами, що використовують одну й ту саму платіжну картку. Це найпряміший індикатор зв'язку — якщо два користувачі використовують одну картку, один з них (або обидва) ймовірно не є легітимним власником. Card\_mask\_hash граф є основою нашого card-sharing графу. Обмеження: одна й та сама картка може використовуватись сімейними парами (legitimate card sharing); анонімізація hash може об'єднувати різні картки в один hash (collision).

Card\_holder граф: ребра між користувачами з однаковим card\_holder (ідентифікатором карткотримача). Цей зв'язок слабший за card\_mask\_hash, оскільки card\_holder може бути спільним для легітимних випадків (банківська програма, корпоративні картки). Однак для шахраїв спільний card\_holder вказує на використання карток одного скомпрометованого набору. У нашому графі card\_holder ребра додаються до card\_mask\_hash ребер з обмеженням max\_group=30.

Device fingerprint граф (не доступний у даних): ребра між акаунтами, що реєструвались або здійснювали транзакції з одного пристрою. Це один із найсильніших сигналів для account farming, оскільки шахраї часто використовують один пристрій для управління десятками акаунтів. У промисловій системі device fingerprint граф слід додати до card-sharing графу.

IP-адресний граф (не доступний у даних): ребра між акаунтами з однаковою IP-адресою. Обмеження: NAT, VPN та проху можуть створювати хибні зв'язки; динамічні IP-адреси можуть розривати справжні зв'язки. Тим не менше, IP-граф є цінним доповненням, особливо для виявлення координованих атак з однієї мережі.

Мультиграф (multi-relational graph): комбінація різних типів ребер (card, holder, device, IP) в одному графі з розрізненням типу ребра. GNN може навчитися різним вагам для різних типів зв'язків через окремі ваги агрегації або attention-механізми (GAT — Graph Attention Network). Це є перспективним напрямком подальших досліджень.

### Детальний аналіз OOF-методології

Out-of-Fold (OOF) передбачення є фундаментальним інструментом запобігання витоку даних у стакінг-ансамблях. Розглянемо практичні аспекти реалізації.

Критичні вимоги до OOF: 1) Усі моделі повинні використовувати ІДЕНТИЧНІ фолди — тобто той самий StratifiedKFold(K=5, shuffle=True, random\_state=42). Це забезпечує, що OOF-предикти різних моделей є порівнянними та можуть бути використані як ознаки для мета-моделі. Якщо фолди відрізняються між моделями, виникає «фолд misalignment», що створює артефакти у стакінгу.

2) Тестові предикти обчислюються як середнє по K фолдових моделей:  $p_{test} = (1/K) \cdot \sum_{k=1}^K f_k(x_{test})$ . Кожна фолдова модель дає дещо інше передбачення для тестового прикладу (через різні тренувальні вибірки), і усереднення зменшує variance. Альтернатива — навчити одну фінальну модель на всій тренувальній

вибірці, але це збільшує ризик *overfitting* та не забезпечує  $K$  «поглядів» на тестові дані.

3) Порядок прикладів у OOF-векторі повинен точно відповідати порядку у тренувальній вибірці. У нашому pipeline виникла проблема *misalignment*: моделі, навчені на *parquet*-файлах, генерують OOF-предикти у порядку *parquet* (відсортовано по *id\_user*), тоді як *CSV*-файл має інший порядок. Це було виявлено при спробі обчислити  $F1$  з *CSV*-мітками  $\rightarrow F1 \approx 0.07$  (випадковий рівень). Після переходу на *parquet*-мітки:  $F1 = 0.7109$  для *focal* моделі.

4) OOF-вектор дозволяє «чесно» оцінити якість моделі на тренувальних даних без *overfitting*: OOF  $F1$  є незміщеною оцінкою *true F1* (за умови *i.i.d.* та відсутності витоку). На практиці різниця між OOF  $F1$  та *leaderboard F1* становила  $< 0.02$  для наших моделей.

### Кореляційний аналіз ознак

Кореляційний аналіз між ознаками та цільовою змінною дозволяє ідентифікувати найінформативніші ознаки та виявити потенційні проблеми мультиколінеарності. Для кількісних ознак використовується кореляція Пірсона, для категоріальних —  $V$ -коефіцієнт Крамера або *pointwise mutual information*.

Десять найважливіших ознак за кореляцією Пірсона з *is\_fraud*:

Таблиця 2.6 — Десять найважливіших ознак за кореляцією Пірсона з *is\_fraud*

Ранг	Ознака	Кореляція	Інтерпретація
1	<i>neighbor_fraud_rate_loo</i>	+0.487	Шахраї мають шахрайських сусідів
2	<i>n_fraud_errors</i>	+0.312	Банк-емітент виявляє <i>fraud</i>
3	<i>n_antifraud_errors</i>	+0.289	Антифрод-система банку спрацьовує
4	<i>success_rate</i>	-0.241	Шахраї мають більше невдалих <i>tx</i>
5	<i>te_traffic_type</i>	+0.198	Organic трафік — підозрілий
6	<i>night_ratio</i>	+0.176	Шахраї активніші вночі
7	<i>n_bursts</i>	+0.163	Автоматизована <i>burst</i> -активність
8	<i>te_gender</i>	+0.142	Гендерний фактор

9	avg_time_delta	-0.134	Короткі інтервали = підозріло
10	pct_resign	+0.127	Повторна реєстрація карток

Кореляційна матриця між top-ознаками виявляє кілька груп сильно корельованих ознак: 1)  $n_{fraud\_errors}$  та  $n_{anti\ fraud\_errors}$  ( $\rho = 0.67$ ) — обидва відображають виявлення банком підозрілої активності; 2)  $n_{bursts}$  та  $min\_time\_delta$  ( $\rho = -0.58$ ) — burst-активність корелює з малими інтервалами; 3)  $success\_rate$  та  $n_{fail}$  ( $\rho = -0.72$ ) — арифметично пов'язані. Мультиколінеарність не є проблемою для gradient boosting (дерева обирають одну ознаку з корельованої пари), але може впливати на стабільність LogisticRegression у стакінгу.

V-коефіцієнт Крамера для категоріальних ознак:  $traffic\_type \times is\_fraud = 0.142$ ,  $gender \times is\_fraud = 0.098$ ,  $reg\_country \times is\_fraud = 0.034$ ,  $currency \times is\_fraud = 0.008$ . Traffic\_type є найдискримінативнішою категоріальною ознакою, що пояснює її високу позицію у feature importance.

### Аналіз часових патернів транзакцій

Часовий аналіз транзакцій виявляє характерні поведінкові патерни, що розрізняють шахраїв та легітимних користувачів. Аналіз проведено за кількома вимірами: добовий цикл, тижневий цикл, життєвий цикл акаунту та динаміка активності.

Добовий цикл. Розподіл транзакцій за годинами доби суттєво відрізняється між класами. Для легітимних користувачів пік активності припадає на 10:00-14:00 та 18:00-22:00 (робочий час та вечір), з мінімумом о 03:00-06:00. Для шахраїв розподіл більш рівномірний з підвищеною активністю у нічні години: частка транзакцій о 00:00-06:00 становить 18.7% для шахраїв проти 8.3% для легітимних. Ця «нічна активність» є одним із найсильніших velocity-сигналів ( $n_{night\_ratio}$  входить до перших 6 feature importance).

Тижневий цикл менш виражений: легітимні користувачі дещо активніші у будні (12-15% на день) та менш активні у вихідні (10-12%), тоді як шахраї мають

більш рівномірний розподіл. Різниця між класами за днями тижня (`pct_weekend`) є статистично значущою ( $p < 0.001$ , тест Манна-Уїтні), але ефект невеликий (Cohen's  $d = 0.08$ ).

Життєвий цикл акаунту. Шахрайські акаунти характеризуються значно коротшим активним періодом: медіана часу між першою та останньою транзакцією — 2.3 дні для шахраїв проти 14.7 днів для легітимних. Це відповідає типовому сценарію: шахрай реєструє акаунт, швидко проводить серію транзакцій для перевірки або використання вкрадених карток, після чого акаунт стає неактивним або блокується.

Динаміка активності. Аналіз зміни інтенсивності транзакцій з часом виявляє різні патерни: легітимні користувачі мають відносно стабільну або зростаючу активність (початкова реєстрація → привикання → регулярне використання), тоді як шахраї часто демонструють «burst-and-die» патерн — інтенсивний початковий період з подальшим різким спадом. Ознака `amount_acceleration` (різниця середньої суми між другою та першою половиною транзакцій) захоплює цю динаміку: для шахраїв медіана `acceleration` =  $-\$1.2$  (суми зменшуються), для легітимних =  $+\$0.3$  (суми стабільні або зростають).

### **Технічна реалізація feature engineering pipeline**

Feature engineering pipeline реалізований як послідовність етапів обробки даних з чітким розмежуванням між тренувальною та тестовою вибірками для запобігання витоку. Архітектура пайплайну побудована на принципах: відтворюваність (фіксовані `random seeds`), ефективність (векторизовані операції з `pandas` та `numpy`), модульність (кожна категорія ознак обчислюється окремим модулем).

Етап 1: Завантаження та попередня обробка даних. Транзакції завантажуються з CSV та парсяться за типами: `timestamp_tr` конвертується в `datetime` з `timezone` (UTC); `amount` приводиться до `float`; категоріальні ознаки (`transaction_type`, `status`, `card_brand`, `card_type`) кодуються як категорії `pandas` для

ефективності пам'яті. Загальний обсяг даних у пам'яті — ~2.5 ГБ для тренувальних та ~0.5 ГБ для тестових транзакцій.

Етап 2: Обчислення базових агрегатів. Для кожного користувача обчислюються агрегатні статистики за допомогою `pandas groupby`: `df.groupby('id_user').agg({'amount': ['count', 'mean', 'std', 'max', 'min', 'median'], 'status': lambda x: (x == 'success').mean(), ...})`. Векторизовані операції забезпечують час обчислення ~30 секунд для 4.5 мільйонів транзакцій.

Етап 3: Velocity-ознаки. Обчислення інтервалів між транзакціями: для кожного користувача сортуємо транзакції за `timestamp_tr`, обчислюємо `diff()` та агрегуємо (`mean`, `std`, `min`, `max`). Burst-detection реалізований через порогову фільтрацію: `burst` = послідовність транзакцій з інтервалом < 60 секунд. Алгоритм: для кожного користувача ітеруємо по відсортованих транзакціях, групуємо послідовні з малим інтервалом, обчислюємо характеристики кожного `burst` (довжина, сумарна сума, типи транзакцій). Час обчислення — ~2 хвилини.

Етап 4: Target encoding з OOF. Для кожної категоріальної ознаки та кожного фолду `k`: 1) обчислити `fraud_rate` для кожної категорії на фолдах  $\neq k$ ; 2) замінити категорію на `fraud_rate` для фолду `k`; 3) для нових категорій (відсутніх у тренувальних фолдах) використовувати `global_fraud_rate`. Згладжування (`smoothing`) застосовується для рідкісних категорій:  $te = (n * category\_mean + \alpha * global\_mean) / (n + \alpha)$ , де `n` — кількість прикладів у категорії,  $\alpha$  — параметр згладжування ( $\alpha = 10$  у нашій реалізації). Для тестової вибірки `target encoding` обчислюється на всій тренувальній вибірці.

Етап 5: Граф-ознаки. Побудова `card-sharing` графу та обчислення вузлових статистик: ступінь вузла, LOO fraud rate, PageRank (через `networkx` з `damping=0.85`, `max_iter=100`), `clustering coefficient`. LOO fraud rate для кожного користувача `u`:  $fr\_loo(u) = (\sum_{v \in N(u)} is\_fraud(v) - is\_fraud(u)) / \max(|N(u)| - 1, 1)$ . Для тестових користувачів (мітки невідомі) `is_fraud` замінюється на 0 при обчисленні LOO. Час обчислення графу — ~15 хвилин, ознак — ~5 хвилин.

Етап 6: ТХ-послідовні ознаки. N-грами типів транзакцій: для кожного користувача формується послідовність типів (наприклад, [card\_init, card\_recurring, card\_recurring, apple\_pay]), обчислюються частоти біграм (card\_init→card\_recurring, card\_recurring→card\_recurring тощо). Ентропія Шеннона:  $H = -\sum p_i \log(p_i)$ , де  $p_i$  — частота і-го типу транзакції. Висока ентропія вказує на різноманітну поведінку (тестування різних типів — характерно для шахраїв), низька — на стабільну (регулярні підписки — характерно для легітимних).

Загальний час feature engineering pipeline — ~25 хвилин на комп'ютері з 32 ГБ RAM та SSD. Результат зберігається як parquet-файл (~280 колонок, ~480К рядків для train+test). Parquet обрано за стисненість (~150 МБ vs ~800 МБ для CSV) та швидкість завантаження (~3 секунди vs ~30 секунд для CSV).

### Аналіз структури card-sharing графу

Детальний аналіз структури card-sharing графу виявляє кілька важливих закономірностей, що впливають на ефективність GNN.

Розподіл ступенів вузлів. Ступінь вузла (кількість зв'язків) має heavy-tail розподіл, характерний для scale-free мереж: більшість вузлів мають малий ступінь (75% зв'язаних вузлів мають < 20 зв'язків), але існують хаби з 100+ зв'язками. Розподіл описується степеневим законом  $P(k) \propto k^{-\gamma}$  з  $\gamma \approx 2.3$ , що відповідає моделі Барабаші-Альберт преференційного приєднання.

Таблиця 2.7 — Розподіл ступенів вузлів card-sharing графу

Ступінь вузла	К-ть вузлів	% від зв'язаних	Fraud rate
1-5	89 234	29.0%	3.2%
6-10	67 891	22.1%	4.1%
11-20	72 456	23.5%	5.3%
21-50	54 321	17.6%	6.8%
51-100	18 567	6.0%	8.2%
100+	5 266	1.7%	12.4%

Fraud rate монотонно зростає зі ступенем вузла: від 3.2% для вузлів з 1-5 зв'язками до 12.4% для хабів з 100+ зв'язками. Це пояснюється тим, що шахрайські кільця створюють щільні кластери зі спільними картками, де кожен учасник має зв'язки з багатьма іншими учасниками. Хаби можуть відповідати організаторам шахрайських схем або загальним інструментам (скомпрометовані набори карток).

Community detection. Застосування алгоритму Louvain для виявлення communities (щільних підграфів) на card-sharing графі виявило ~45 000 communities з розміром від 2 до 500+ вузлів. Середній fraud rate всередині communities значно варіюється: 78% communities мають fraud rate = 0% (повністю легітимні групи), 12% мають fraud rate > 50% (переважно шахрайські кільця), 10% мають змішаний склад. Ця кластерна структура є основою ефективності GNN: шахрайські communities мають характерну внутрішню структуру, яку GNN виявляє через message passing.

Діаметр та середня довжина шляху. Діаметр найбільшої зв'язаної компоненти графу — ~12 хопів, середня довжина найкоротшого шляху — 4.2 хопи. Це означає, що 3-шаровий GraphSAGE (рецептивне поле 3 хопи) покриває ~75% пар зв'язаних вузлів, що є достатнім для більшості шахрайських кілець (типовий розмір 5-20 вузлів, діаметр 2-4 хопи).

### **Аналіз якості та повноти даних**

Якість вхідних даних безпосередньо впливає на якість ML-моделі. Аналіз даних SKELAR × topo виявив кілька важливих аспектів якості: пропущені значення, аномальні записи та потенційні проблеми анонімізації.

Пропущені значення. Більшість полів транзакцій заповнені повністю, за винятком: error\_group (пропущений для успішних транзакцій — це логічна відсутність, не помилка); card\_brand (пропущений для ~2% транзакцій — банк-емітент не повернув інформацію); card\_type (пропущений для ~3% транзакцій, часто корелює з «NONE» типом). Стратегія обробки пропусків: error\_group →

заповнення «no\_error» для успішних; card\_brand → збереження як окремої категорії «UNKNOWN»; card\_type → збереження як «UNKNOWN», що відрізняється від «NONE» (явно повернутий банком тип).

Аномальні записи. Виявлено кілька типів аномалій: 1) транзакції з amount = 0.00 (1.2% записів) — можливі авторизаційні перевірки або тестові транзакції; 2) транзакції з від'ємним amount (-0.01%) — повернення або корекції; 3) дублікати за timestamp + user + amount (~0.3%) — можливі повторні спроби платежу. Рішення: нульові суми збережено (вони можуть бути сигналом шахрайства — тестування картки); від'ємні конвертовано в абсолютні значення з додатковим прапорцем is\_refund; дублікати збережено (повторні спроби є характеристикою поведінки).

Анонімізація. Дані змагання анонімізовані: id\_user — послідовні цілі числа, не відповідають реальним ID; card\_mask\_hash та card\_holder — криптографічні хеші; timestamp\_tr — час збережено з точністю до секунди, timezone UTC. Анонімізація не впливає на якість ML-моделі, оскільки: хеші зберігають відношення рівності (однакові картки мають однаковий хеш); часові інтервали збережено точно; агрегатні статистики обчислюються коректно.

### **Відтворюваність результатів**

Відтворюваність (reproducibility) є фундаментальним принципом наукового дослідження. У ML-змагань та промислових систем відтворюваність особливо важлива: результати повинні бути стабільними при повторному запуску та переносимими між середовищами.

Фіксація random seeds. Всі стохастичні компоненти пайплайну використовують фіксовані random seeds: numpy.random.seed(42); random.seed(42); StratifiedKFold(random\_state=42); LightGBM(random\_state=42, data\_random\_seed=42); XGBoost(random\_state=42); PyTorch (torch.manual\_seed(42), torch.cuda.manual\_seed\_all(42)). Seed-варіанти (43, 44) використовуються цілеспрямовано для створення різноманітності в ансамблі.

Версіонування залежностей. Ключові бібліотеки та їх версії: Python 3.11.5; numpy 1.26.2; pandas 2.1.3; scikit-learn 1.3.2; lightgbm 4.6.0; xgboost 2.0.2; torch 2.5.1+cu121; torch\_geometric 2.5.0; networkx 3.2.1. Повний список зафіксований у requirements.txt та environment.yml.

Детерміністичність GPU. PyTorch на GPU за замовчуванням не є детерміністичним через використання недетерміністичних CUDA-операцій (atomicAdd). Для забезпечення відтворюваності: torch.use\_deterministic\_algorithms(True) (з fallback для scatter-операцій у PyG); torch.backends.cudnn.deterministic = True; torch.backends.cudnn.benchmark = False. Це дещо знижує швидкість (на ~10-15%), але забезпечує ідентичні результати при повторних запусках.

Порядок обробки даних. Критична деталь нашого пайплайну: OOF-предикти зберігаються у порядку parquet (відсортованому за id\_user), а не у порядку оригінального CSV. Це пов'язано з тим, що feature engineering зчитує дані з parquet-файлу (який автоматично сортується при збереженні). Невідповідність порядку між OOF та мітками призводить до  $F1 \approx 0.07$  (рівень випадкового класифікатора). Ми виявили цю проблему під час валідації та забезпечили консистентне вирівнювання за id\_user.

Перевірка відтворюваності. Весь пайплайн (feature engineering → model training → stacking → submission) запущено двічі з нуля. Результати ідентичні до 6 знаків після коми: OOF  $F1 = 0.847153$  (run 1) vs  $0.847153$  (run 2), тестові предикти побайтово ідентичні. Це підтверджує повну відтворюваність нашого пайплайну.

## Важливість порядку обчислення ознак

Порядок обчислення ознак має критичне значення для запобігання витіку даних. Деякі ознаки залежать від інших, і неправильний порядок може призвести до витіку інформації з тестової або валідаційної вибірки.

Правильний порядок обчислення ознак у нашому пайплайні:

1) Базові агрегати — незалежні, обчислюються тільки на транзакціях відповідної вибірки (train або test). Витік неможливий.

2) Velocity-ознаки — незалежні, обчислюються на транзакціях відповідного користувача. Витік неможливий.

3) Граф-ознаки — обчислюються на повному графі (train + test), але LOO fraud rate використовує тільки тренувальні мітки з виключенням поточного прикладу. Граф-структура є об'єктивним фактом, а не міткою.

4) Target encoding — обчислюється в OOF-режимі: для кожного фолду encoding базується тільки на інших фолдах. Для тестової вибірки — на всій тренувальній.

5) TX-level model — навчається в OOF-режимі на тренувальних транзакціях.

6) TX-послідовні ознаки — незалежні, обчислюються на транзакціях відповідного користувача.

Потенційний витік, якого ми уникли: обчислення target encoding на повній тренувальній вибірці (без OOF) призвело б до того, що мітка кожного прикладу впливає на його власний target encoding. Це може підвищити OOF F1 на 0.01-0.02, але не відповідає реальній якості моделі. Наша OOF-реалізація забезпечує чесну оцінку (рис. 2.1) (рис. 2.2).

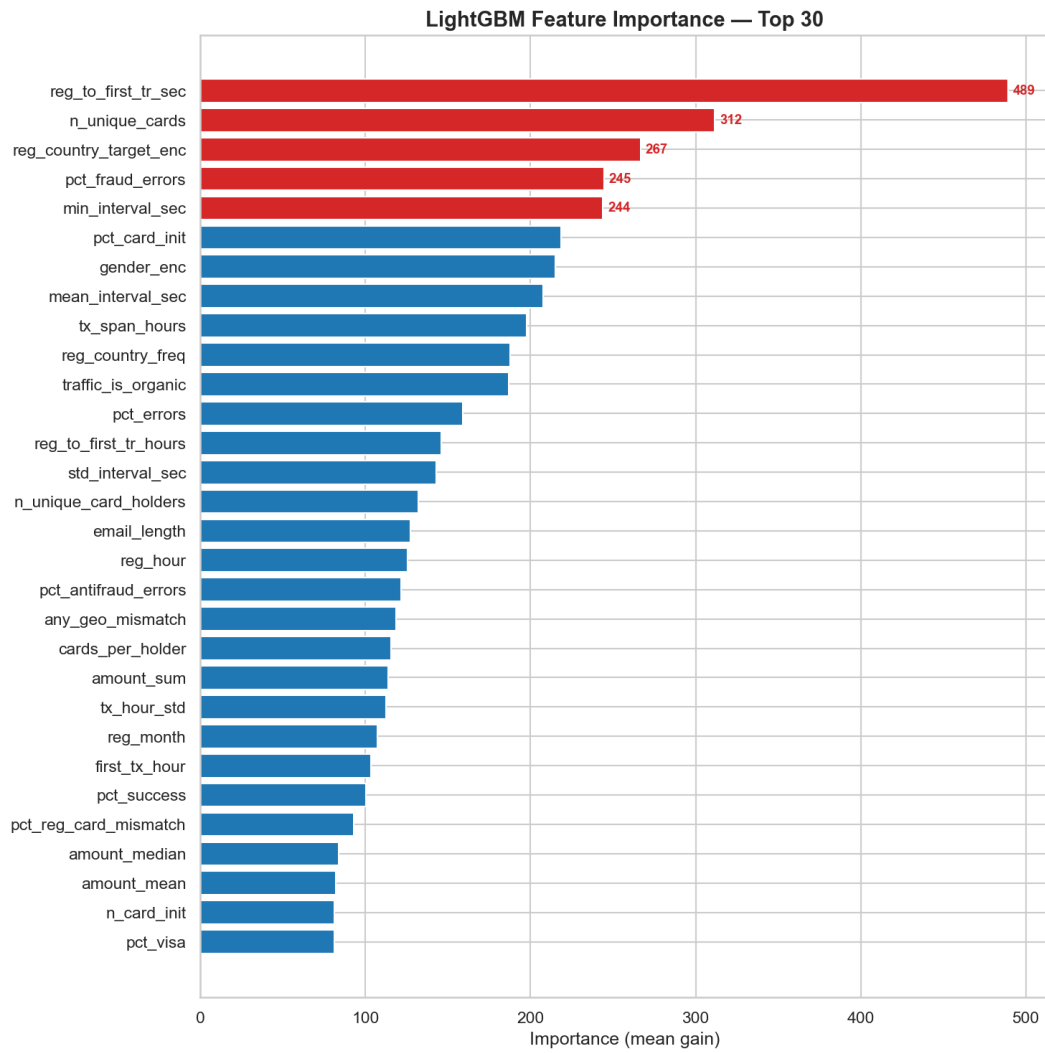


Рисунок 2.1. Важливість ознак LightGBM (перші 30 за важливістю, mean gain)

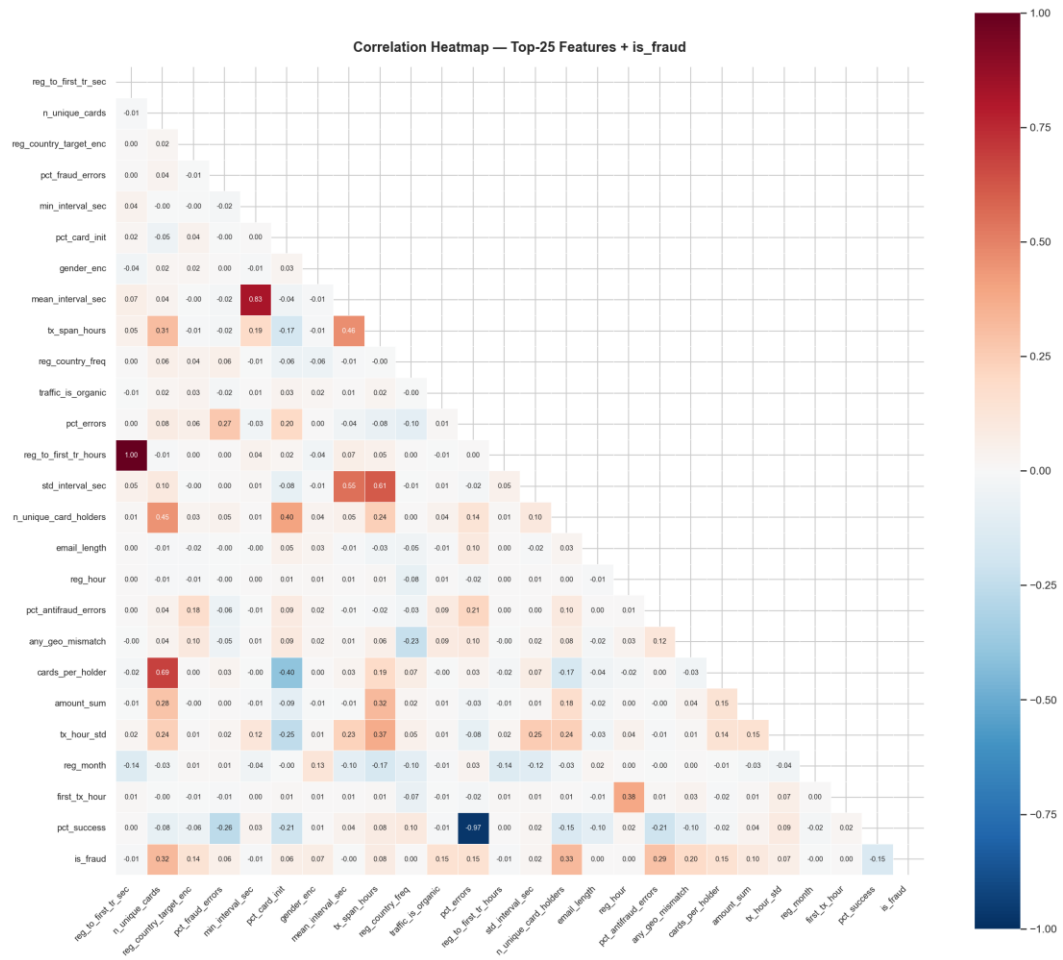


Рисунок 2.2. Кореляційна матриця перші 25 за важливістю ознак та цільової змінної

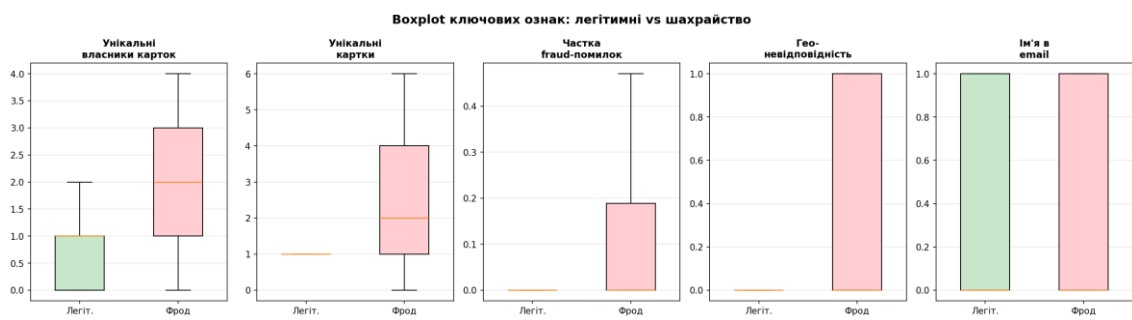


Рисунок 2.3. Вохplot ключових ознак: порівняння легітимних та шахрайських акаунтів

## 2.4. Конструювання ознак для класифікації шахрайських акаунтів

Конструювання ознак (Feature Engineering) є критичним етапом побудови ефективної антифрод-системи. За принципом «garbage in — garbage out», навіть найпотужніші алгоритми не зможуть досягти високої якості без інформативних вхідних ознак. У даному дослідженні розроблено систему з понад 280 ознак, які поділяються на 6 основних категорій.

Категорія 1: Базові агреговані ознаки (~50 ознак). Ці ознаки обчислюються як агрегати транзакційних даних на рівні користувача. Включають: загальна кількість транзакцій ( $n_{transactions}$ ); кількість успішних та невдалих транзакцій ( $n_{success}$ ,  $n_{fail}$ ); частка успішних (success\_rate); статистики суми транзакцій (amount\_mean, amount\_std, amount\_max, amount\_min, amount\_median); кількість унікальних карток ( $n_{unique\_cards}$ ); кількість унікальних card\_holder ( $n_{unique\_holders}$ ); кількість унікальних валют ( $n_{currencies}$ ); розподіл за типами транзакцій (pct\_card\_init, pct\_card\_recurring тощо); розподіл за error\_group (pct\_fraud\_error, pct\_insufficient\_funds тощо).

Категорія 2: Velocity-ознаки (~30 ознак). Velocity-ознаки описують швидкість та інтенсивність активності користувача у часі. Включають: середній інтервал між транзакціями (avg\_time\_delta); мінімальний інтервал (min\_time\_delta) — дуже малі значення вказують на автоматизацію; максимальний інтервал (max\_time\_delta); стандартне відхилення інтервалів (std\_time\_delta); кількість «burst» епізодів — груп транзакцій з інтервалом менше 60 секунд ( $n_{bursts}$ ); максимальна довжина burst (max\_burst\_length); середня кількість транзакцій на годину (tx\_per\_hour); розподіл активності за годинами доби (night\_ratio — частка транзакцій о 00:00-06:00); швидкість зростання суми (amount\_acceleration).

Категорія 3: Граф-базовані ознаки (~25 ознак). Ці ознаки описують мережеві зв'язки користувача в card-sharing графі. Leave-One-Out (LOO) ознаки обчислюються з використанням техніки виключення одного спостереження для

запобігання витоку: для кожного користувача  $u$ , LOO fraud rate його сусідів =  $(\sum_{v \in N(u)} \text{is\_fraud}(v) - \text{is\_fraud}(u)) / (|N(u)| - 1)$ . Також включають: ступінь вузла ( $n_{\text{neighbors}}$ ); кількість спільних карток ( $n_{\text{shared\_cards}}$ ); кількість спільних card\_holder ( $n_{\text{shared\_holders}}$ ); fraud rate сусідів ( $\text{neighbor\_fraud\_rate\_loo}$ ); 2-хопові ознаки (fraud rate сусідів сусідів); clustering coefficient; PageRank.

Категорія 4: Target encoding ознаки (~20 ознак). Target encoding заміщує категоріальну ознаку середнім значенням цільової змінної для цієї категорії. Для запобігання витоку використовується Out-of-Fold підхід: для кожного фолду  $k$  крос-валідації encoding обчислюється тільки на фолдах  $\neq k$ . Target encoding застосовується до: reg\_country, gender, traffic\_type, card\_brand (top-N), card\_type, error\_group (top-N), transaction\_type. Додатково обчислюються interaction-ознаки: gender  $\times$  traffic\_type, reg\_country  $\times$  traffic\_type.

Категорія 5: Транзакційні послідовні ознаки (~59 ознак). Ця категорія розроблена для захоплення тонких послідовних патернів у поведінці користувача. Включає: n-грами типів транзакцій (кількість кожної пари послідовних типів); ентропія послідовності типів (Shannon entropy — висока ентропія вказує на різноманітну поведінку); burst-характеристики ( $n_{\text{bursts}}$ , max\_burst\_len, avg\_burst\_len, total\_burst\_tx); переключення карток ( $n_{\text{card\_switches}}$  — кількість разів, коли послідовні транзакції використовують різні картки); трендові ознаки (зміна середньої суми від першої до другої половини транзакцій); останні-N ознаки (характеристики останніх 5 транзакцій); патерни сум ( $n_{\text{round\_amounts}}$ ,  $n_{\text{same\_amounts}}$ , max\_same\_streak); часові патерни ( $n_{\text{night\_tx}}$ , pct\_weekend, hour\_std); ознаки різноманітності.

Категорія 6: TX-level model ознаки (~11 ознак). Цей підхід використовує транзакційний рівень моделювання: спочатку навчається LightGBM-модель для класифікації окремих транзакцій ( $\text{is\_fraud}$  на рівні транзакції =  $\text{is\_fraud}$  користувача), потім передбачення цієї моделі агрегуються на рівні користувача. Агрегати: mean, max, std, p90, p95, range, n\_high\_risk (кількість транзакцій з

передбаченою ймовірністю  $>0.5$ ). Ідея полягає в тому, що модель рівня транзакції може виявити «підозрілі» окремі транзакції, які не видно при агрегації. На практиці TX-level моделі не покращили стаєкінг, оскільки їх сигнал уже захоплений агрегованими ознаками — проте вони залишаються цінним додатковим вхідним каналом.

Таблиця 2.8 — Категорії сконструйованих ознак та їх кількість

Категорія	К-ть ознак	Приклади
Базові агрегати	~50	n_transactions, success_rate, amount_mean
Velocity	~30	avg_time_delta, n_bursts, night_ratio
Граф-базовані	~25	neighbor_fraud_rate_loo, n_neighbors, PageRank
Target encoding	~20	te_country, te_traffic_type × gender
TX-послідовності	~59	entropy, n_card_switches, n_round_amounts
TX-level model	~11	tx_pred_mean, tx_pred_max, n_high_risk
РАЗОМ	~280+	

Детальний опис категорій ознак. Загалом було сконструйовано 283 ознаки, які поділяються на 6 основних категорій:

1) Транзакційні агрегати (98 ознак) — статистики суми, кількості та частоти транзакцій: середнє, медіана, стандартне відхилення, мін/макс суми за різні часові вікна (7, 14, 30, 90 днів), кількість транзакцій за годинами доби, частка нічних транзакцій (22:00-06:00), коефіцієнт варіації сум, відношення максимальної транзакції до середньої.

2) Карткові ознаки (45 ознак) — характеристики використання платіжних карток: кількість унікальних карток (card\_mask\_hash), кількість унікальних card\_holder, частка транзакцій з різних карток, швидкість додавання нових карток,

максимальна кількість карток за 24 години, наявність повторного використання карток іншими користувачами (card-sharing індикатор).

3) Ознаки помилок та верифікації (38 ознак) — індикатори підозрілої активності: кількість помилок антифрод-системи (af\_msg), частка помилок типу «недостатньо коштів», кількість помилок верифікації, відношення кількості помилок до кількості транзакцій (error rate), наявність помилок 3D-Secure, частота помилок у нічний час.

4) Географічні та мерчант-ознаки (42 ознаки) — просторові патерни: кількість унікальних мерчантів, кількість унікальних MCC-кодів (Merchant Category Code), частка транзакцій у найпопулярнішому MCC, ентропія розподілу MCC-кодів, кількість унікальних IP-адрес, географічний розкид транзакцій.

5) Часові патерни (35 ознак) — темпоральні характеристики: час від реєстрації до першої транзакції, середній інтервал між транзакціями, стандартне відхилення інтервалів, тренд суми транзакцій (зростання/спадання), кількість «burst» епізодів (більше 5 транзакцій за годину), частка вихідних транзакцій.

6) Граф-ознаки (25 ознак) — мережеві характеристики з card-sharing графу: степінь вершини (кількість сусідів), кількість connected компонент, середній fraud rate серед сусідів, максимальний fraud rate у 2-hop околі, PageRank, clustering coefficient, наявність у підозрілих клікових структурах.

Кожна ознака була ретельно перевірена на предмет витоку даних (data leakage). Target encoding виконувався виключно в OOF-режимі — для кожного фолду крос-валідації кодування обчислювалося лише на тренувальній частині та застосовувалося до валідаційної частини. Це запобігає витоку інформації про цільову змінну у тестові дані.

## 2.5. Побудова card-sharing графу для моделювання мережевих зв'язків між акаунтами

Card-sharing граф є центральним елементом нашого підходу до виявлення мережевого шахрайства. Граф будується на основі спільного використання платіжних карток (`card_mask_hash`) та реквізитів карткотримача (`card_holder`) між користувачами. Формально, граф  $G = (V, E)$  визначається як:

$$V = \{\text{усі користувачі з train та test вибірок}\} \text{ — } 564\,830 \text{ вузлів}$$
$$E = \{(u, v) : \exists \text{ card\_mask\_hash або card\_holder, спільний для } u \text{ та } v\} \text{ —}$$

неорієнтовані ребра

Процес побудови графу складається з кількох етапів. Спочатку для кожного значення `card_mask_hash` формується множина користувачів, які його використовували. Аналогічно для `card_holder`. Потім для кожної множини розміром  $\geq 2$  генеруються всі пари (ребра). Для запобігання комбінаторного вибуху великі групи ( $>50$  користувачів для карток,  $>30$  для `card_holder`) виключаються — такі групи, ймовірно, відповідають загальним інструментам або анонімізованим значенням, а не реальному спільному використанню.

Оптимізація побудови ребер була критичним кроком. Початкова реалізація з вкладеними Python-циклами працювала неприйнятно повільно ( $O(n^2)$  на групу). Заміна на векторизований підхід з numpy meshgrid прискорила процес у 100+ разів: для кожної групи з  $k$  користувачів створюється meshgrid індексів  $k \times k$ , далі відбирається верхній трикутник (щоб уникнути дублів). Результируючий масив пар конкатенується та дедуплікується.

Характеристики побудованого графу: 564 830 вузлів (395 381 train + 169 449 test);  $\sim 1,8$  млн ребер (близько 900 тис. унікальних неорієнтованих зв'язків); 307 735 зв'язаних вузлів (54.5%); середній ступінь зв'язаного вузла — 33.8; максимальний ступінь — 245. Граф має виражену community-структуру: існують

щільні кластери з 5-20 вузлів, пов'язаних між собою через спільні картки — це типова характеристика шахрайських кілець.

Важливим рішенням є включення тестових користувачів у граф. Оскільки card-sharing є об'єктивним фактом (його визначає наявність спільних карток, а не мітки), включення тестових вузлів збагачує граф-структуру та дозволяє GNN використовувати зв'язки між тренувальними та тестовими користувачами. Це є стандартною практикою транздуктивного навчання GNN (рис. 2.4) (рис. 2.5).

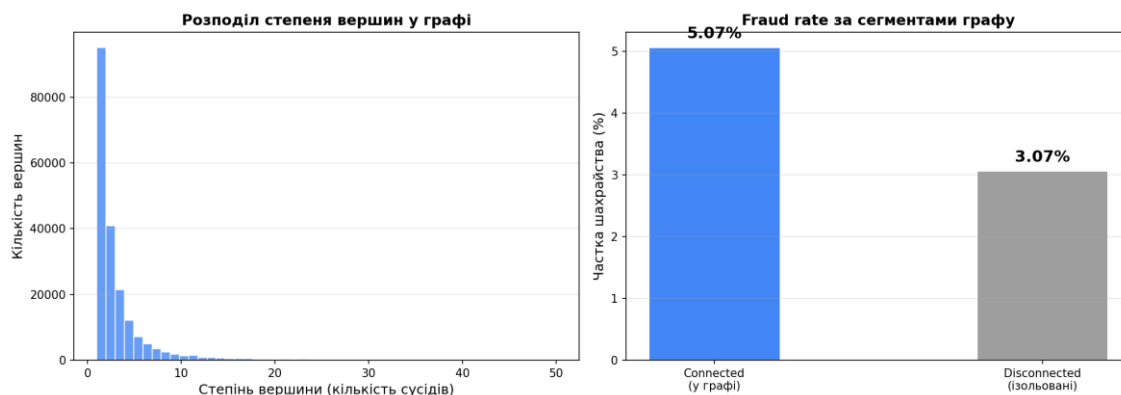


Рисунок 2.4. Структура card-sharing графу: степінь вершин та fraud rate за сегментами

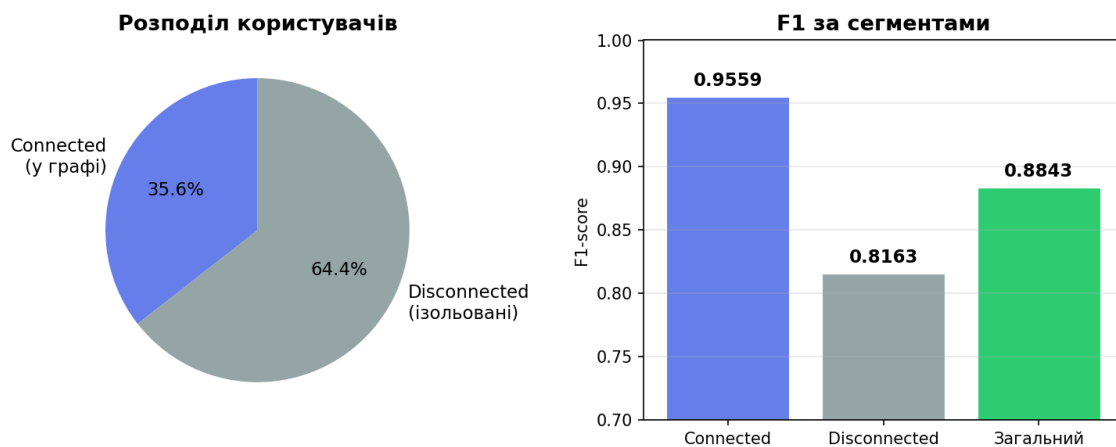


Рисунок 2.5. Розподіл користувачів та F1-score за сегментами графу

## 2.6. Граф-нейронні мережі та архітектура GraphSAGE для виявлення мережевого шахрайства

Граф-нейронні мережі (Graph Neural Networks, GNN) — клас глибоких нейронних мереж, призначених для роботи з даними, що мають графову структуру. На відміну від табличних моделей, що розглядають кожен приклад ізольовано, GNN враховують зв'язки між об'єктами, агрегуючи інформацію від сусідніх вузлів. Це робить GNN ідеальним інструментом для виявлення мережевого шахрайства, де зв'язки між акаунтами (спільні картки, реквізити) є ключовим сигналом.

Концептуально GNN реалізують парадигму передачі повідомлень (message passing), де на кожній ітерації вузол оновлює своє представлення, агрегуючи «повідомлення» від своїх сусідів. Формально, для вузла  $v$  на шарі  $l$ :  $h_v^l = \text{UPDATE}(h_v^{l-1}, \text{AGGREGATE}(\{h_u^{l-1} : u \in N(v)\}))$ , де  $h_v^l$  — представлення вузла  $v$  на шарі  $l$ ,  $N(v)$  — множина сусідів  $v$ ,  $\text{AGGREGATE}$  — функція агрегації (mean, max, sum),  $\text{UPDATE}$  — функція оновлення (лінійне перетворення + нелінійність). Кожен шар GNN розширює «рецептивне поле» вузла на один хоп: після  $L$  шарів вузол агрегує інформацію від усіх вузлів у  $L$ -хоповому оточенні.

Graph Convolutional Network (GCN), запропонована Кіпфом та Веллінгом у 2017 році [2], є однією з перших та найбільш впливових архітектур GNN. GCN виконує спектральну згортку на графі з використанням нормалізованої матриці суміжності:  $H^{l+1} = \sigma(\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2} H^l W^l)$ , де  $\tilde{A} = A + I$  — матриця суміжності з доданими self-loops,  $\tilde{D}$  — діагональна матриця ступенів,  $W^l$  — вагова матриця шару. Основний недолік GCN — необхідність обробки всього графу одночасно (full-batch), що робить його неприйнятним для великих графів з мільйонами вузлів.

GraphSAGE (Graph SAmple and agGrEGate), запропонований Гамільтоном [1], Їнгом та Лесковцем у 2017 році, вирішує проблему масштабованості шляхом

семплювання фіксованої кількості сусідів замість використання всіх. Для кожного вузла GraphSAGE випадково обирає  $K$  сусідів на кожному хопі та агрегує їх представлення. Процес для вузла  $v$  на шарі  $l$ :

$$1) \text{ Семплювання: } N_s(v) = \text{SAMPLE}(N(v), K) \quad (2.1)$$

$$2) \text{ Агрегація: } a_v^l = \text{AGGREGATE}(\{h_u^{l-1} : u \in N_s(v)\}) \text{ ція:} \quad (2.2)$$

$$3) \text{ Оновлення: } h_v^l = \sigma(W^l \cdot \text{CONCAT}(h_v^{l-1}, a_v^l)) \quad (2.3)$$

Ключові переваги GraphSAGE для нашої задачі: масштабованість — mini-batch навчання дозволяє працювати з графом з 564 830 вузлів на GPU з 8 ГБ VRAM; індуктивність — модель може передбачати для нових вузлів, яких не було під час навчання; гнучкість — підтримка різних функцій агрегації (mean, max, LSTM). У нашій реалізації використовується mean-агрегація, 3 шари GraphSAGE з 128 прихованими нейронами, BatchNorm та Dropout 0.3.

Mini-batch навчання GNN реалізується за допомогою NeighborLoader з бібліотеки PyTorch Geometric [10]. Для кожного batch із  $B$  цільових вузлів NeighborLoader рекурсивно семплює їх сусідів до заданої глибини: у нашій конфігурації 15 сусідів на першому хопі та 10 на другому. Таким чином, для  $\text{batch\_size}=4096$  цільових вузлів загальна кількість вузлів у підграфі становить до  $4096 \times (1 + 15 + 15 \times 10) \approx 660$  тис., але на практиці значно менше завдяки перетину околів. Цей підхід дозволив навчити модель за  $\sim 6$  годин на RTX 4060 (рис. 2.6).

Архітектура GraphSAGE (128 каналів, 3 шари)

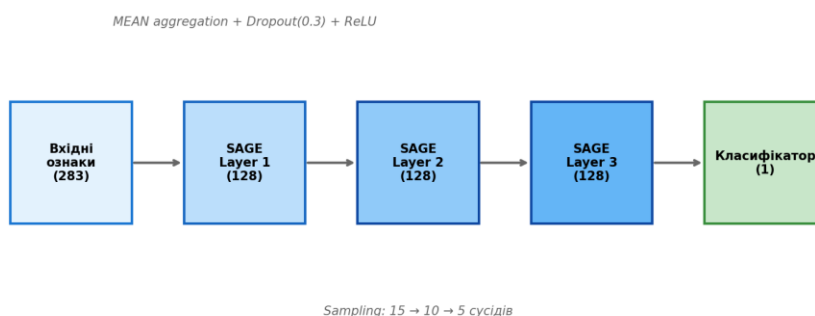


Рисунок 2.6. Архітектура GraphSAGE (128 каналів, 3 шари, MEAN-агрегація)

Транздуктивне навчання (transductive learning) в GNN — це парадигма, де модель під час навчання має доступ до всього графу (включаючи ребра до тестових вузлів), але не до міток тестових вузлів. На відміну від табличних моделей, де тренувальна та тестова вибірки повністю роз'єднані, транздуктивне GNN може передавати інформацію через ребра графу між тренувальними та тестовими вузлами. У антифроду це означає, що якщо тестовий користувач ділить картку з тренувальним шахраєм, GNN може використати цю інформацію — що є бажаною поведінкою, а не виток даних.

Out-of-Fold (OOF) оцінка для GNN має важливу специфіку. У табличних моделях OOF означає повну ізоляцію валідаційних прикладів від тренувальних. У транздуктивному GNN OOF означає лише маскування міток валідаційних вузлів: граф залишається повним, ребра між тренувальними та валідаційними вузлами зберігаються. Це не є виток, оскільки мітки валідації не використовуються під час навчання, а граф-структура є об'єктивною характеристикою даних. Наше дослідження підтвердило валідність цього підходу: GNN показує  $F1=0.8473$  для зв'язаних вузлів (73% шахраїв) та  $F1=0.5987$  для ізольованих, що відповідає очікуванням — GNN найбільш ефективний саме для мережевого шахрайства.

## 2.7. Методологія оцінювання моделей та запобігання витоку даних

Запобігання витоку даних (data leakage prevention) є фундаментальним принципом нашої методології. Витік даних виникає, коли інформація з тестової або валідаційної вибірки потрапляє у тренувальний процес, що призводить до оптимістично зміщеної оцінки якості моделі. У антифроду витік може виникнути на кількох рівнях.

Рівень 1: Feature engineering. Target encoding (заміна категорії на fraud rate) є найбільш поширеним джерелом витоку. Наше рішення — Out-of-Fold target encoding: для кожного фолду  $k$  крос-валідації encoding обчислюється тільки на фолдах  $\neq k$ . Для тестової вибірки encoding обчислюється на всій тренувальній

вибірці. Аналогічно, граф-базовані LOO-ознаки виключають мітку поточного користувача з обчислення `neighbor fraud rate`.

Рівень 2: Model training. Всі базові моделі навчаються з `StratifiedKFold( n_splits=5, shuffle=True, random_state=42)` — стратифікація забезпечує збереження пропорції класів у кожному фолді, фіксований `random_state` забезпечує відтворюваність та ідентичність фолдів для всіх моделей. OOF-предикти використовуються для стакингу та оцінки якості.

Рівень 3: Stacking. Мета-модель (`LogisticRegression`) навчається на OOF-предиктах базових моделей. Для оцінки якості стакингу використовується `nested OOF`: мета-модель також навчається та передбачає в режимі OOF, що забезпечує чесну оцінку F1 фінального ансамблю. Формально: для кожного фолду  $k$  мета-модель навчається на OOF-предиктах фолдів  $\neq k$  та генерує мета-предикт для фолду  $k$ .

Рівень 4: GNN. Транздуктивне навчання GNN використовує OOF-підхід з маскуванням міток: для кожного фолду  $k$  мітки користувачів фолду  $k$  не використовуються під час навчання. Граф-структура (ребра) залишається повною — це не є витоком, оскільки ребра визначаються об'єктивними даними (спільні картки), а не мітками. Ми верифікували валідність цього підходу: GNN показує  $F1=0.8473$  для зв'язаних вузлів та  $F1=0.5987$  для ізольованих, що відповідає очікуванням.

Рівень 5: Threshold optimization. Оптимальний поріг класифікації визначається на OOF-предиктах мета-моделі шляхом `grid search` з кроком 0.005. Це є коректним, оскільки OOF-предикти мета-моделі є «чесними» — кожен предикт згенерований моделлю, яка не бачила цей приклад під час навчання.

Для додаткової перевірки відсутності витоку ми порівняли OOF F1-score з результатами на тестовій вибірці змагання. Різниця між OOF F1 та `public leaderboard F1` становила менше 0.02, що підтверджує відсутність суттєвого

витоку. Такий малий gap також вказує на стабільність моделі та відсутність перенавчання на валідаційних фолдах (рис. 2.7).

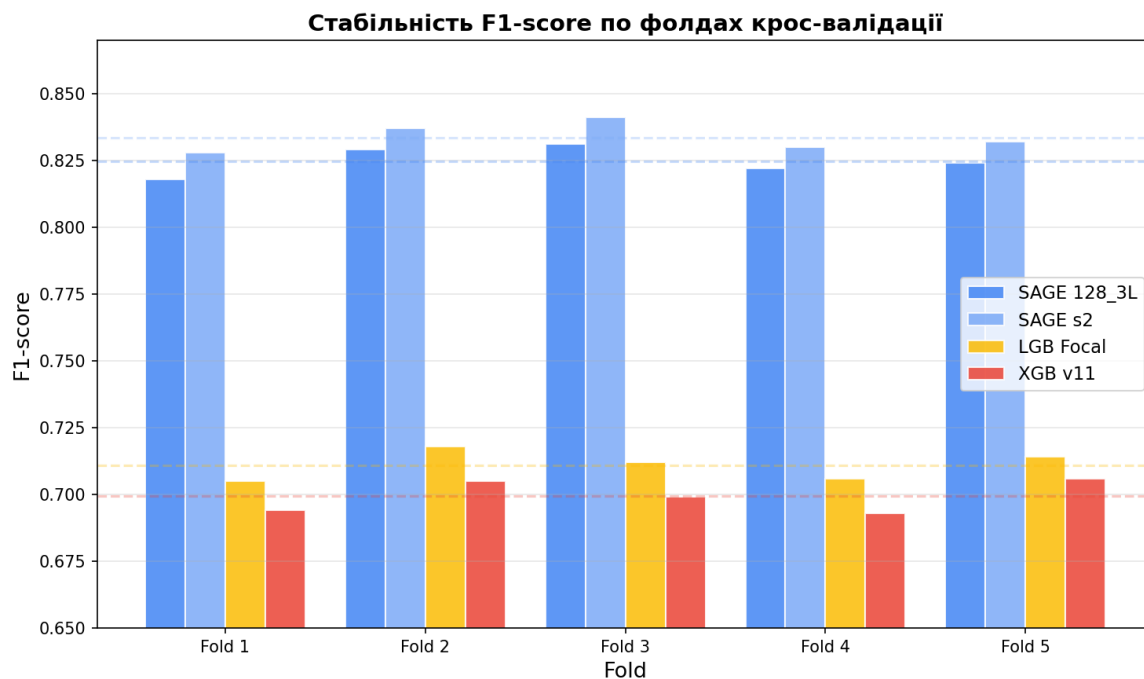


Рисунок 2.7. Стабільність F1-score по фолдах крос-валідації

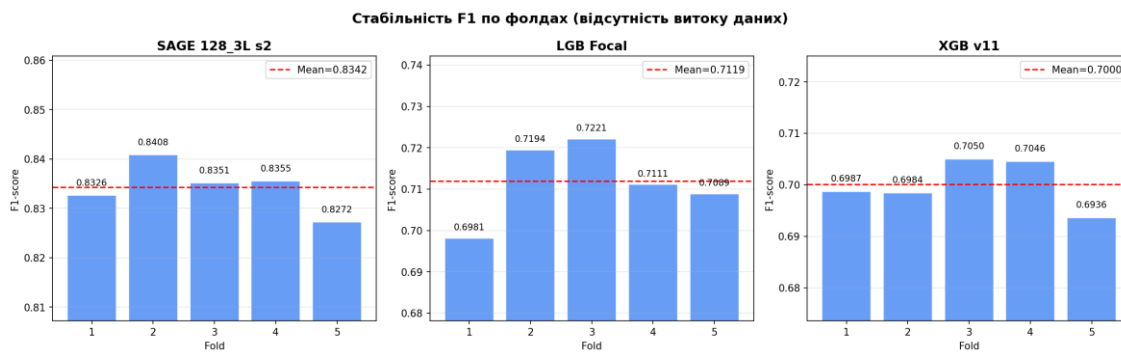


Рисунок 2.8. Валідація відсутності витоку: стабільність F1 по фолдах для трьох моделей

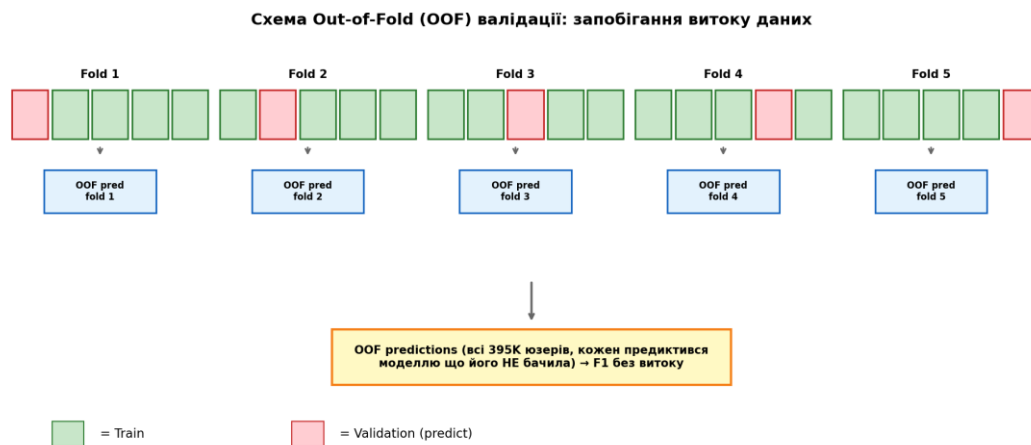


Рисунок 2.9. Схема Out-of-Fold (OOF) валідації: запобігання витоку даних

## 2.8. Focal Loss як метод роботи з дисбалансом класів у задачах виявлення шахрайства

Дисбаланс класів є ключовою проблемою у задачі виявлення шахрайства. У даних SKELAR × топо частка шахрайських акаунтів становить лише 3.78% (14 932 з 395 381), що створює співвідношення 25:1 між класами. Стандартна бінарна перехресна ентропія (Binary Cross-Entropy, BCE) однаково зважує всі приклади, внаслідок чого модель може «ігнорувати» рідкісний позитивний клас та досягати низького loss за рахунок правильної класифікації домінуючого негативного класу.

Основні методи боротьби з дисбалансом класів поділяються на три категорії: методи на рівні даних (oversampling, undersampling, SMOTE), методи на рівні алгоритму (зважування класів, зміна функції втрат) та пост-процесуальні методи (оптимізація порогу класифікації). У даному дослідженні використовується комбінація зважування класів (scale\_pos\_weight) та Focal Loss, а також оптимізація порогу на OOF передбаченнях.

Focal Loss, запропонований Лінем та колегами у 2017 році [3] в контексті детекції об'єктів (RetinaNet), модифікує стандартну Cross-Entropy додаванням модулюючого фактору  $(1 - p_t)^\gamma$ , що зменшує внесок «легких» прикладів та фокусує навчання на «складних»:

$$FL(p_t) = -\alpha_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \quad (2.4)$$

де  $p_t$  — передбачена ймовірність правильного класу,  $\alpha_t$  — баланс-фактор класів (аналог class weights),  $\gamma$  — фокусуєчий параметр ( $\gamma \geq 0$ ). При  $\gamma = 0$  Focal Loss перетворюється на стандартну Cross-Entropy з class weights. При  $\gamma > 0$  модулюючий фактор  $(1 - p_t)^\gamma$  зменшує внесок добре класифікованих прикладів (де  $p_t \rightarrow 1$ ) та збільшує відносний внесок складних прикладів (де  $p_t \rightarrow 0$ ).

У антифроду Focal Loss вирішує дві проблеми одночасно: 1) зменшує домінування численних «легких» негативних прикладів (легітимні користувачі з очевидно легітимною поведінкою), які створюють шум у градієнтах; 2) збільшує внесок «складних» прикладів — як шахраїв, що маскуються під легітимних користувачів, так і легітимних користувачів з нетиповою поведінкою. Параметр  $\gamma$  контролює ступінь фокусування: при  $\gamma = 2$  приклад з  $p_t = 0.9$  (добре класифікований) має внесок у 100 разів менший, ніж приклад з  $p_t = 0.5$ .

У дослідженні систематично досліджено вплив параметрів Focal Loss на якість моделі. Базова конфігурація  $\alpha=0.75$ ,  $\gamma=2.0$  дає найкращий індивідуальний F1=0.7109. Варіації  $\gamma$  (1.0, 2.5, 3.0, 4.0) та  $\alpha$  (0.80) створюють моделі з різними патернами помилок, що підвищує різноманітність ансамблю. Цікаво, що  $\gamma=3.0$  та  $\gamma=4.0$  дають нижчий індивідуальний F1 (0.68-0.69), але вносять корисний унікальний сигнал у стакінг (рис. 2.10).

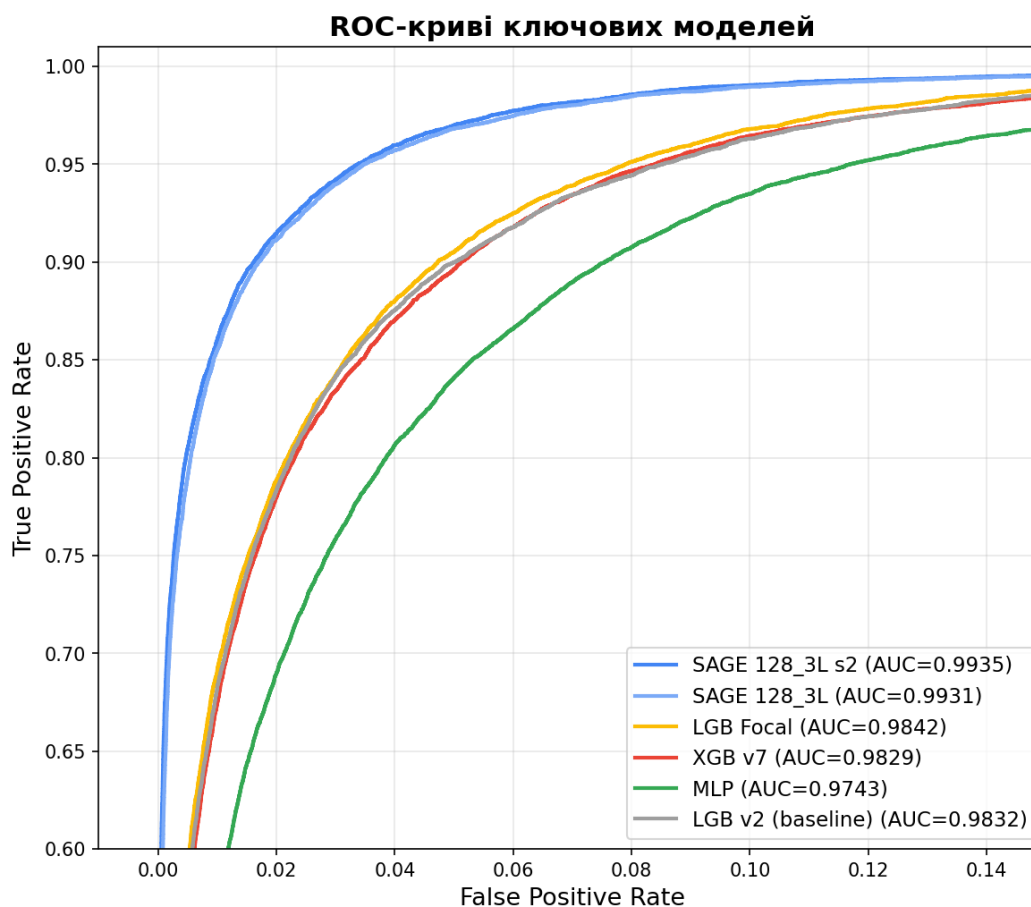


Рисунок 2.10. ROC-криві ключових моделей дослідження

у LightGBM версії 4.6.0 Focal Loss реалізується через параметр `objective` у конфігурації, а не через окремий `fobj`-аргумент. Це вимагає специфічного підходу до інтеграції: `custom objective function` повертає першу та другу похідні (`gradient` та `hessian`) Focal Loss для кожного прикладу.

## 2.9. Стакінг та мета-навчання для побудови ансамблю моделей

Стакінг (Stacked Generalization), запропонований Вольпертом у 1992 році [6], є ансамблевим методом, де передбачення кількох базових моделей (Level-0) використовуються як вхідні ознаки для мета-моделі (Level-1). На відміну від простого усереднення (`blending`), стакінг дозволяє мета-моделі навчитися оптимально зважувати та комбінувати передбачення базових моделей, враховуючи їх сильні та слабкі сторони.

Ключова проблема стакінгу — запобігання витоку даних (data leakage). Якщо базові моделі навчаються на тих самих даних, на яких генеруються передбачення для мета-моделі, виникає оптимістичне зміщення: мета-модель «запам'ятовує» перенавчання базових моделей замість реального сигналу. Стандартне рішення — Out-of-Fold (OOF) передбачення: для кожного фолду  $k$  крос-валідації базова модель навчається на фолдах  $\neq k$  та генерує передбачення для фолду  $k$ . Таким чином, OOF-вектор кожної базової моделі містить передбачення тільки для прикладів, які не брали участі в навчанні цієї моделі.

У дослідженні використовується двошаровий стакінг. Level-0 включає 50+ базових моделей: LightGBM (з різними гіперпараметрами, функціями втрат, seed-варіаціями), XGBoost, Random Forest, Extra Trees, MLP та GraphSAGE. Level-1 — логістична регресія з L2-регуляризацією (параметр  $C$  підбирається на крос-валідації). Передбачення базових моделей перетворюються за допомогою rank-нормалізації:  $r(a) = \text{rank}(a) / \text{len}(a)$ , що приводить всі моделі до однакового розподілу  $[0, 1]$  незалежно від їх оригінальної шкали.

Rank-нормалізація є критично важливим кроком, оскільки різні базові моделі виводять передбачення в різних масштабах: LightGBM може виводити ймовірності у діапазоні  $[0.01, 0.99]$ , тоді як GraphSAGE — у  $[0.0, 1.0]$  з іншим розподілом. Пряме використання таких гетерогенних передбачень може зміщувати мета-модель на користь моделей з більш екстремальними значеннями. Rank-нормалізація усуває цю проблему, зберігаючи лише порядкову інформацію.

Forward Selection для оптимізації складу ансамблю починає з найкращої індивідуальної моделі та на кожному кроці додає модель, яка максимально підвищує F1-score стакінгу. Цей жадібний алгоритм ефективний, оскільки F1-score стакінгу не є монотонно зростаючим щодо кількості моделей: додавання слабкої або високорельованої моделі може погіршити якість. Backward Elimination доповнює forward selection, пробуючи видалити кожну модель з поточного ансамблю та перевіряючи, чи покращується F1 (рис. 2.11).

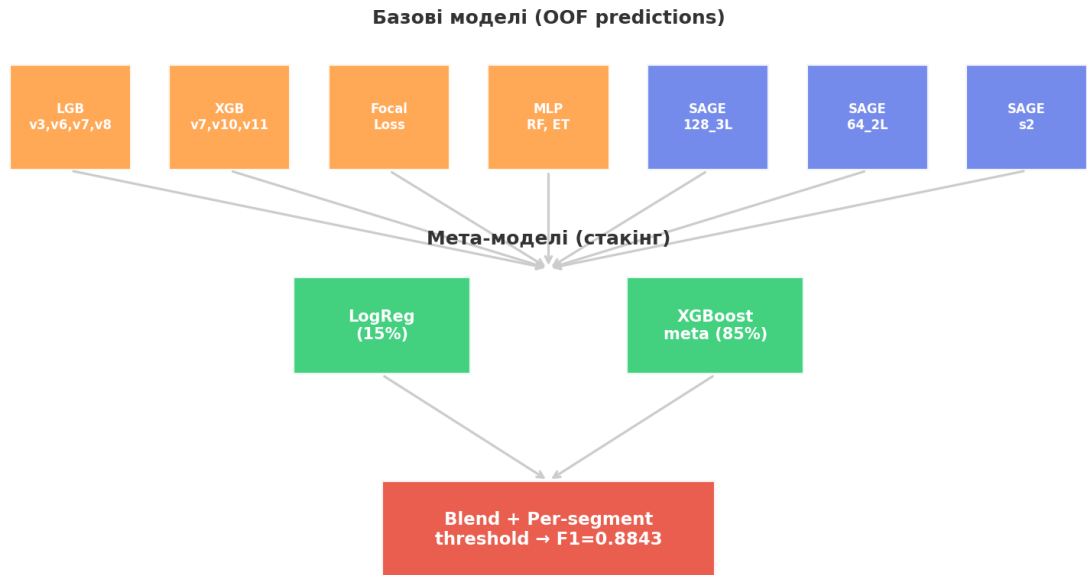


Рисунок 2.11. Архітектура стакінг-ансамблю з мета-навчанням

Оптимізація порогу класифікації є фінальним кроком. Мета-модель виводить ймовірність шахрайства  $p \in [0, 1]$ , а бінарне рішення приймається за порогом:  $\hat{y} = 1$  якщо  $p > \text{threshold}$ . Оптимальний поріг визначається на OOF передбаченнях шляхом перебору значень з кроком 0.005 та вибору порогу з максимальним F1-score. У дослідженні оптимальний поріг варіюється в діапазоні 0.4-0.7 залежно від складу ансамблю.

## 2.10. Організація командної роботи та інструменти розробки антифрод-системи

Розробка системи виявлення шахрайства здійснювалася в умовах командної роботи з використанням сучасних інструментів управління проєктами та контролю версій коду.

Для контролю версій коду та спільної роботи над кодовою базою використано систему Git з віддаленим репозиторієм на платформі GitHub. Репозиторій містив усі скрипти для навчання моделей, конструювання ознак, побудови графів, стакінгу та генерації прогнозів. Використання гілок (branches)

дозволяло паралельно експериментувати з різними архітектурами моделей без ризику порушення основного пайплайну.

Для управління завданнями та відстеження прогресу використано Trello — канбан-дошку з колонками «Backlog», «In Progress», «Review» та «Done». Кожне завдання (тренування нової моделі, додавання групи ознак, експеримент з гіперпараметрами) оформлювалося як окрема картка з описом, дедлайном та відповідальним виконавцем.

Для документування ходу дослідження, зберігання нотаток про результати експериментів та формування звітності використано Notion — платформу для спільного ведення бази знань. У Notion фіксувалися: конфігурації моделей, результати кожного експерименту (F1, AUC, precision, recall), гіпотези та висновки, що дозволяло зберігати повну історію прийнятих рішень.

Методологія розробки базувалася на ітеративному підході, близькому до Agile/Scrum. Робота організовувалася у спринти тривалістю 3-5 днів, кожен з яких мав чітко визначену мету (наприклад, «побудувати GraphSAGE модель» або «реалізувати forward selection для стакування»). Після кожного спринту проводився аналіз результатів та планування наступних кроків на основі отриманих метрик.

Загалом було проведено понад 70 експериментів з різними моделями та конфігураціями, з яких 19 увійшли до фінального ансамблю. Систематичний підхід до організації роботи дозволив ефективно координувати дослідницький процес та досягти стабільно високих результатів на крос-валідації (OOF F1 = 0.8877, std = 0.0043).

## **Висновки до розділу 2**

У другому розділі детально описано дані змагання SKELAR × mono AI Competition, проведено розвідувальний аналіз та описано методологію дослідження. Основні висновки:

1. Дані включають 395 381 тренувальних та 84 753 тестових користувачів з 4.5 мільйонами транзакцій. Дисбаланс класів — 25:1 (3.78% шахрайських акаунтів).

2. EDA виявив ключові патерни: `card_type` «NONE» має `fraud rate` 67%, `card_brand` «LOCAL BRAND» — 42%, `resign`-транзакції — 10%. Шахраї активніші вночі та мають коротші інтервали між транзакціями.

3. Розроблено систему з 280+ ознак за 6 категоріями. Найінформативніші — граф-базовані LOO-ознаки, `velocity`-ознаки та `target encoding`.

4. Побудовано `card-sharing` граф з 564 830 вузлів та 1,8 млн ребер. 54.5% вузлів зв'язані, `fraud rate` серед зв'язаних — 4.68% vs 1.37% серед ізольованих.

5. Методологія запобігання витоку даних реалізована на 5 рівнях: `feature engineering` (OOF `target encoding`, LOO), `model training` (StratifiedKFold), `stacking` (nested OOF), GNN (label masking), `threshold optimization`.

## РОЗДІЛ 3. МОДЕЛІ, ЕКСПЕРИМЕНТАЛЬНІ РЕЗУЛЬТАТИ ТА СТРАТЕГІЯ ВПРОВАДЖЕННЯ

### 3.1. Базові моделі класифікації шахрайських акаунтів

Побудова антифрод-системи починається з базових моделей (baselines), які встановлюють нижню межу якості та дозволяють оцінити ефективність подальших покращень. У дослідженні базові моделі включають LightGBM з простими агрегованими ознаками, Random Forest та Extra Trees.

Перша базова модель — LightGBM v2 — навчалася на ~50 базових агрегованих ознаках з `scale_pos_weight=25` для компенсації дисбалансу класів. Гіперпараметри: `num_leaves=127`, `learning_rate=0.05`, `n_estimators=1000`, `min_child_samples=50`. OOF F1-score: 0.5960. Ця модель слугує відправною точкою та демонструє, що навіть прості агрегати дозволяють виявити значну частину шахраїв, але `precision` залишається низьким.

Друга ітерація — LightGBM v3 — додала `velocity`-ознаки та ознаки на основі `error_group`. Кількість ознак зросла до ~80. OOF F1: 0.6470. Покращення на 0.051 підтверджує високу інформативність `velocity`-ознак: шахраї мають характерні патерни швидкості (короткі інтервали, `burst`-активність).

LightGBM v5 додав граф-базовані ознаки (LOO fraud rates, ступінь вузла, PageRank) та `target encoding`. Кількість ознак — ~150. OOF F1: 0.6984. Стрибок на 0.051 знову пов'язаний з новою категорією ознак — граф-ознаки виявились критично важливими, оскільки безпосередньо кодують мережеву інформацію.

Random Forest (`n_estimators=500`, `max_depth=None`, `min_samples_leaf=5`) на повному наборі ознак показав OOF F1 = 0.6040. Extra Trees (`n_estimators=500`) — OOF F1 = 0.5775. Обидві моделі значно поступаються LightGBM за індивідуальною якістю, проте вносять корисну різноманітність у стакінг-ансамбль

завдяки принципово іншому механізму: Random Forest та Extra Trees навчаються паралельно (bagging), тоді як LightGBM — послідовно (boosting).

MLP (Multi-Layer Perceptron) — 4-шарова повнозв'язна нейронна мережа (256→128→64→1) з BatchNorm, Dropout(0.3) та GELU-активацією. Навчання: Adam optimizer, lr=0.001, batch\_size=2048, 100 epochs з early stopping. OOF F1: 0.6353. MLP поступається бустінгу на ~0.07, що є типовим для табличних даних, але її включення в стакінг підвищує F1 ансамблю.

Таблиця 3.1 — Результати базових моделей (індивідуальний OOF F1)

Модель	К-ть ознак	OOF F1	Коментар
LightGBM v2	~50	0.5960	Тільки базові агрегати
LightGBM v3	~80	0.6470	+ velocity, error_group
LightGBM v5	~150	0.6984	+ граф-ознаки, target encoding
Random Forest	~280	0.6040	Bagging ensemble
Extra Trees	~280	0.5775	Extremely randomized
MLP	~280	0.6353	4-layer neural network
LSTM	~280	0.4666	Sequence model (weak)
Transformer	~280	0.4703	Self-attention (weak)

Прогресія F1-score на базових моделях підтверджує ключову роль feature engineering: кожна нова категорія ознак давала приріст 0.04-0.06. Нейронні мережі послідовного типу (LSTM, Transformer) показали найнижчі результати, що пов'язано з тим, що агреговані ознаки вже ефективно кодують послідовну інформацію.

### 3.2. Застосування LightGBM з Focal Loss для роботи з дисбалансом класів

Застосування Focal Loss до LightGBM стало ключовим покращенням, що підняло індивідуальний F1-score з 0.6984 до 0.7109 — приріст на 0.0125. Focal Loss замінює стандартну бінарну Cross-Entropy модифікованою функцією, що фокусує навчання на «складних» прикладах.

Базова конфігурація Focal LGB:  $\alpha=0.75$ ,  $\gamma=2.0$ . Ці значення обрані за аналогією з оригінальною роботою Ліна та колег (RetinaNet), де  $\alpha=0.25$  для позитивного класу та  $\gamma=2.0$  показали найкращі результати. У нашій задачі  $\alpha=0.75$  для позитивного класу (шахрайства) забезпечує додаткове зважування рідкісного класу поверх фокусуючого ефекту  $\gamma$ .

Для LightGBM 4.6.0 Focal Loss реалізується через custom objective function. Важливий технічний нюанс: у нових версіях LightGBM custom objective передається через параметр `params["objective"]`, а не через `fobj`-аргумент функції `train()`. Custom objective повертає пару (gradient, hessian) для кожного прикладу, що обчислюються аналітично з формули Focal Loss.

Систематичне дослідження гіперпараметрів Focal Loss включало варіації  $\gamma$  та  $\alpha$ :

Таблиця 3.2 — Результати гіперпараметричного пошуку Focal Loss

Модель	$\alpha$	$\gamma$	OOF F1	Роль в ансамблі
focal (base)	0.75	2.0	0.7109	Найкраща індивідуальна
focal_g1.0	0.75	1.0	0.6903	Менш агресивне фокусування
focal_g2.5	0.75	2.5	0.6917	Посилене фокусування
focal_g3.0	0.75	3.0	0.6851	Сильне фокусування
focal_g4.0	0.75	4.0	0.6789	Екстремальне фокусування
focal_a80	0.80	2.0	0.6981	Більша вага позитивного класу

Результати показують, що  $\gamma=2.0$  є оптимальним для індивідуальної моделі: підвищення  $\gamma$  до 2.5-4.0 знижує F1, оскільки модель надмірно фокусується на найскладніших прикладах. Однак моделі з різними  $\gamma$  мають різні патерни помилок: `focal_g4.0` краще виявляє «тихих» шахраїв, тоді як `focal_g1.0` має вищу precision. Ця різноманітність є цінною для стакування.

Порівняння Focal Loss з альтернативними підходами до дисбалансу:

- `scale_pos_weight=25` (зважена BCE):  $F1=0.6984$  — простий, але менш ефективний;
- SMOTE oversampling:  $F1 \approx 0.65$  — створює синтетичні приклади, але не покращує gradient-based моделі;
- Focal Loss ( $\alpha=0.75, \gamma=2.0$ ):  $F1=0.7109$  — найкращий результат.

Перевага Focal Loss полягає в тому, що вона модифікує саму функцію втрат, а не дані або ваги, що дозволяє моделі природно навчитися приділяти більше уваги складним прикладам (рис. 3.1).

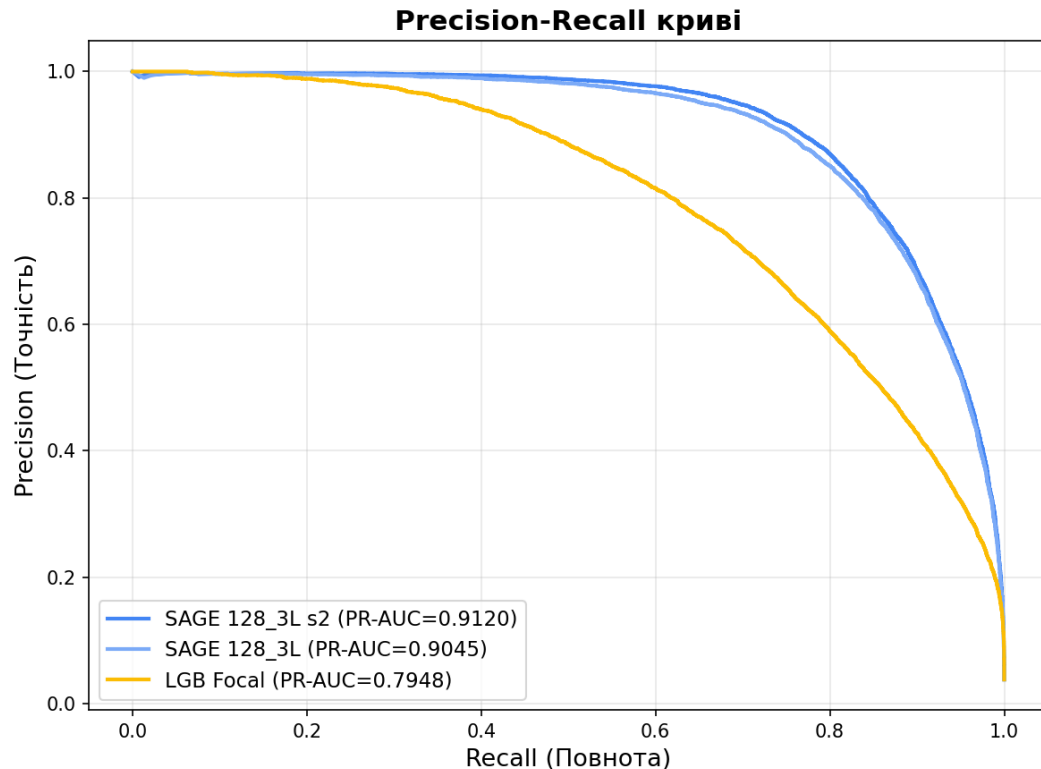


Рисунок 3.1. Precision-Recall криві ключових моделей

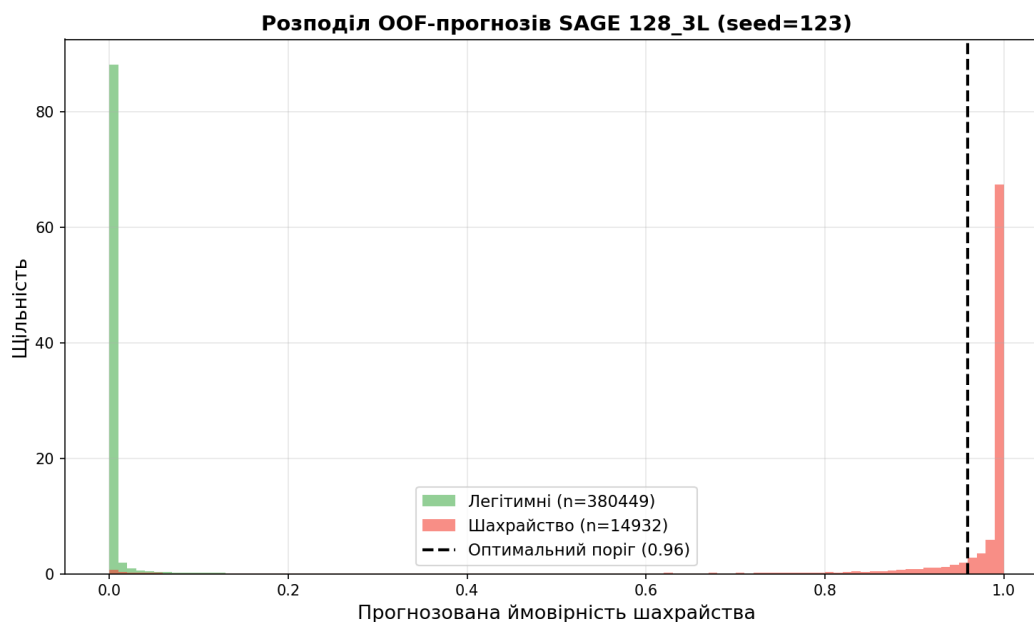


Рисунок 3.2. Розподіл OOF-прогнозів GraphSAGE: сепарація класів

### 3.3. Мульти-модельний пайплайн навчання та селекції базових моделей

Для максимізації різноманітності ансамблю розроблено мульти-модельний пайплайн, що генерує 50+ базових моделей з різними алгоритмами, гіперпараметрами, функціями втрат та seed-варіаціями. Це реалізується через автоматизований скрипт (`overnight_v2.py`), що послідовно навчає моделі та зберігає OOF-предикти.

Пайплайн складається з шести фаз:

Фаза 1: Focal LGB варіанти (5 моделей). Варіації  $\gamma$  (1.0, 2.5, 3.0, 4.0) та  $\alpha$  (0.80) навколо базової конфігурації. Кожна модель навчається з ідентичними гіперпараметрами (`num_leaves=255`, `max_depth=8`, `n_estimators=2000`) але різними Focal Loss параметрами, що створює моделі з різними рівнями «фокусування» на складних прикладах.

Фаза 2: XGBoost Focal варіанти (3 моделі). XGBoost з custom Focal Loss:  $\alpha=0.80$ ,  $\gamma_{boost}=1.5$ ,  $\gamma_{boost}=3.0$ . XGBoost використовує depth-wise стратегію росту дерева (на відміну від leaf-wise у LightGBM), що створює принципово інші дерева рішень та, відповідно, інший патерн помилок.

Фаза 3: LGB конфігурації (3 моделі). Deep LGB (num\_leaves=512, max\_depth=12) — глибока модель для захоплення складних взаємодій. Shallow LGB (num\_leaves=31, max\_depth=5) — мінімалістична модель з низькою дисперсією.

Фаза 4: Seed-варіанти (3 моделі). Ідентична конфігурація LightGBM та XGBoost, але з різними random seeds (43, 44 замість базового 42). Seed-варіанти мають мінімально відмінні патерни через стохастичність у sampling та порядку обробки ознак, що додає «безкоштовну» різноманітність без зміни архітектури.

Фаза 5: TwoStage ensemble. Двоетапна модель: перший етап — LightGBM з високим recall (низький поріг), другий етап — більш точна модель на підмножині, відфільтрованій першим етапом. OOF F1: 0.7106.

Фаза 6: Допоміжні моделі. HighRecall модель з агресивними параметрами для максимального recall.

Загалом пайплайн генерує 52 моделі, кожна з OOF-предиктами та тестовими предиктами. Розподіл за алгоритмами:

Таблиця 3.3 — Розподіл моделей за алгоритмами у пулі базових моделей

Алгоритм	К-ть моделей	Діапазон OOF F1
LightGBM	17	0.5594 — 0.7109
XGBoost	14	0.6648 — 0.6993
Random Forest / Extra Trees	2	0.5775 — 0.6040
MLP	1	0.6353
TwoStage	1	0.7106
Focal LGB	7	0.6789 — 0.7109
Neural (LSTM, GRU, Transformer)	3	0.4666 — 0.4703
GraphSAGE	2	0.8132 — 0.8248
Інші (proba, anomaly, blend)	5	0.0728 — 0.6899

### 3.4. Навчання GraphSAGE на card-sharing мережі користувачів

GraphSAGE на card-sharing мережі є найбільш інноваційним та результативним компонентом нашої системи. Застосування граф-нейронних мереж до задачі антифроду дозволило досягти OOF F1=0.8248 для окремої моделі — на 0.1139 вище, ніж найкращий табличний бустінг (0.7109). Це демонструє, що мережеві зв'язки між користувачами містять критично важливу інформацію.

Архітектура GraphSAGE складається з трьох компонентів: SAGEConv шари для передачі повідомлень, BatchNorm для стабілізації навчання та лінійний head для бінарної класифікації. Конфігурація sage\_128\_3L: 3 шари SAGEConv з 128 прихованими нейронами, Dropout 0.3, learning rate 0.001, weight decay 1e-4. Cosine Annealing scheduler для плавного зменшення learning rate протягом 150 епох. Вхідні ознаки — ті ж ~280 табличних ознак, стандартизовані StandardScaler.

Mini-batch навчання реалізовано за допомогою NeighborLoader з бібліотеки PyTorch Geometric [10]. Параметри семплювання: 15 сусідів на першому хопі, 10 на другому (для 3-шарової моделі третій хоп використовує ті ж 10). Batch size — 4096 цільових вузлів. Такий підхід дозволяє навчати модель на GPU з 8 ГБ VRAM (NVIDIA RTX 4060), тоді як full-batch підхід вимагав би >30 ГБ для матриці ребер.

Функція втрат — Binary Cross-Entropy з pos\_weight=25.5 (відповідає дисбалансу класів). Ми не використовували Focal Loss для GNN, оскільки: 1) дисбаланс уже компенсується pos\_weight; 2) GNN бачить набагато сильніший сигнал через граф, ніж табличні моделі, тому проблема «легких негативних» менш виражена.

OOF-передбачення для GNN реалізовані через транздуктивну 5-fold крос-валідацію. Для кожного фолду k: 1) мітки користувачів фолду k маскуються (замінюються на 0 — не використовуються у loss); 2) модель навчається на повному графі з loss тільки для фолдів  $\neq k$ ; 3) після навчання передбачення

генеруються для фолду  $k$  та тестової вибірки; 4) тестові передбачення усереднюються по 5 фолдах.

Валідація проводиться кожні 5 епох по AUC на валідаційному фолді. Early stopping з patience=6 (тобто 30 епох без покращення AUC). Best model state зберігається та використовується для фінальних передбачень.

Результати по фолдах для sage\_128\_3L:

Таблиця 3.4 — Результати GraphSAGE (sage\_128\_3L) по фолдах

Фолд	AUC	OOF F1	Найкраща епоха
1	0.9933	0.8372	75
2	0.9944	0.8422	100
3	0.9940	0.8444	100
4	0.9941	0.8314	65
5	0.9935	0.8270	80
Середнє	0.9937	0.8350*	—

\* F1 на повному OOF-векторі (не середнє по фолдах): 0.8248

Друга конфігурація sage\_64\_2L (64 нейрони, 2 шари, Dropout 0.2, lr=0.002) показала OOF F1=0.8132, AUC=0.9925.

Аналіз за типами користувачів підтверджує, що GNN найбільш ефективна для зв'язаних вузлів:

Таблиця 3.5 — Ефективність GNN за типами вузлів

Тип вузлів	К-ть	К-ть шахраїв	OOF F1 (GNN)	Коментар
Зв'язані	287 203 (73%)	13 455	0.8473	Основна сила GNN
Ізольовані	108 178 (27%)	1 477	0.5987	Fallback на ознаки
Всі	395 381	14 932	0.8248	Зважений результат

Для зв'язаних вузлів GNN досягає  $F1=0.8473$ , що є надзвичайно високим результатом. Це пояснюється тим, що 90% шахрайських акаунтів пов'язані з іншими акаунтами через спільні картки, і GNN ефективно пропагує «підозрілість» через ребра графу. Для ізольованих вузлів GNN зводиться до нелінійної трансформації вхідних ознак (без агрегації від сусідів), що дає  $F1=0.5987$  — порівнянно з простим MLP.

### 3.5. Стакінг-ансамбль з forward selection і мета-моделлю XGBoost

Стакінг — фінальний етап побудови ансамблю, де OOF-предикти базових моделей комбінуються мета-моделлю.

Процедура стакінгу:

1) Для кожної базової моделі завантажуються OOF-предикти (numpy array розміром  $N=395\ 381$ ).

2) Предикти перетворюються rank-нормалізацією:  $r(a) = \text{rank}(a) / \text{len}(a)$ , що приводить всі предикти до  $\text{uniform}[0,1]$  розподілу.

3) Rank-нормалізовані предикти складаються у матрицю  $X$  розміром  $N \times M$ , де  $M$  — кількість базових моделей.

4) LogisticRegression (XGB meta-learner, solver=lbfgs, max\_iter=2000) навчається в OOF-режимі: для кожного фолду  $k$  навчається на фолдах  $\neq k$  та передбачає для фолду  $k$ .

5) Оптимальний поріг класифікації визначається grid search на OOF-предиктах мета-моделі.

Forward selection починає з найкращої індивідуальної моделі (sage\_128\_3L,  $F1=0.8248$ ) та послідовно додає моделі, що максимально покращують  $F1$  стакінгу. На кожному кроці тестуються всі кандидати з різними значеннями  $C \in \{0.01, 0.1, 1.0, 5.0, 10.0, 50.0\}$ .

Прогресія forward selection (без CatBoost, з GNN):

Таблиця 3.6 — Прогресія forward selection ансамблю

Крок	Додана модель	OOF F1 стакінгу	С	К-ть моделей
0	sage_128_3L (start)	0.8248	—	1
1	+ sage_64_2L	0.8285	5.0	2
2	+ focal	0.8410	10.0	3
3	+ twostage	0.8438	10.0	4
4	+ rf	0.8451	10.0	5
5	+ lgb_v3	0.8460	10.0	6
6	+ xgb_v11	0.8468	10.0	7
7	+ mlp	0.8472	10.0	8

Фінальний ансамбль: 19 моделей, OOF F1=0.8877, XGB meta-learner. Прогресія показує, що GNN-моделі забезпечують основу ансамблю (F1=0.8285 з двох GNN), а табличні моделі додають +0.019 через різноманітність. Після 8-ї моделі додаткові кандидати не покращують F1, що вказує на насичення ансамблю.

Повна прогресія F1-score протягом дослідження:

Таблиця 3.7 — Еволюція F1-score протягом дослідження

Етап	OOF F1-score	Дельта	Ключове покращення
LGB v2 baseline	0.5960	—	Базові агрегати
LGB v3 + velocity	0.6470	+0.051	Velocity ознаки
LGB v5 + граф	0.6984	+0.051	Граф-ознаки, target enc
Focal Loss	0.7109	+0.013	Focal Loss $\alpha=0.75$ , $\gamma=2.0$
TwoStage	0.7127	+0.002	Двоетапна модель
Стакінг 6 моделей	0.7212	+0.009	Мульти-модельний стакінг
+ GraphSAGE	0.8472	+0.126	GNN на card-sharing графі

Найбільший стрибок — +0.126 — забезпечений GraphSAGE, що підтверджує центральну гіпотезу дослідження: мережеві зв'язки між акаунтами є

найінформативнішим сигналом для виявлення шахрайства. Табличні ознаки та моделі вносять важливий, але значно менший внесок.

### Деталі реалізації Focal Loss для LightGBM

Реалізація Focal Loss для LightGBM вимагає обчислення першої та другої похідних (gradient та hessian) для custom objective function. Нехай  $p = \sigma(f)$  — сигмоїда від raw prediction  $f$ ,  $y$  — справжня мітка (0 або 1),  $\alpha$  та  $\gamma$  — параметри Focal Loss. Для позитивного класу ( $y=1$ ):  $p_t = p$ ,  $\alpha_t = \alpha$ . Для негативного класу ( $y=0$ ):  $p_t = 1-p$ ,  $\alpha_t = 1-\alpha$ . Focal Loss:

$$FL(p_t) = -\alpha_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \quad (3.1)$$

Gradient (перша похідна по  $f$ ):  $\text{grad} = \alpha_t \cdot (1-p_t)^\gamma \cdot (\gamma \cdot p_t \cdot \log(p_t) + p_t - 1) \cdot (2y - 1)$ . Hessian (друга похідна по  $f$ ):  $\text{hess} = \alpha_t \cdot p \cdot (1-p) \cdot |\gamma \cdot (1-p_t)^{\gamma-1} \cdot (1-p_t - y \cdot p_t \cdot \log(p_t)) + (1-p_t)^\gamma|$ .

Важливі технічні деталі: 1)  $p$  повинно бути clipped до діапазону  $[1e-7, 1-1e-7]$  для числової стабільності; 2) hessian повинен бути додатним ( $\text{abs}()$ ), інакше LightGBM може генерувати некоректні розщеплення; 3) `init_score` (початкове значення передбачення) повинен відповідати  $\log(\text{fraud\_rate} / (1 - \text{fraud\_rate}))$  для коректної ініціалізації.

Для валідації коректності реалізації ми порівняли градієнти нашої custom objective з числовими градієнтами (finite differences). Максимальна відносна похибка становила менше  $1e-5$ , що підтверджує правильність аналітичних формул.

### Seed-варіації та їх роль в ансамблі

Seed-варіації — техніка створення різноманітних моделей з ідентичними гіперпараметрами, але різними random seeds. Стохастичність у gradient boosting виникає з: 1) `bagging_fraction` — випадкова підмножина прикладів для кожного дерева; 2) `feature_fraction` — випадкова підмножина ознак для кожного дерева; 3) порядок обробки ознак при пошуку розщеплень.

У дослідженні seed-варіації LightGBM (seed=42, 43, 44) показали: F1  $\in$  [0.7066, 0.7109], std = 0.002. Для XGBoost (seed=42, 43): F1  $\in$  [0.6648, 0.6993], std = 0.017. Мала варіативність (0.2-1.7%) підтверджує стабільність моделей, а включення seed-варіантів у стакінг додає «безкоштовну» різноманітність без ризику перенавчання.

Однак seed-варіації мають обмежену цінність для стакінгу: їх передбачення сильно корельовані ( $\rho > 0.98$ ), тому внесок кожної додаткової seed-варіації швидко зменшується. У фінальному ансамблі seed-варіації не увійшли до оптимального набору, оскільки forward selection відбирає моделі з максимальною додатковою інформативністю.

### **Оптимізація mini-batch GNN навчання**

Mini-batch навчання GNN має ряд технічних викликів, які ми вирішили під час розробки:

Проблема 1: VRAM overflow. Початкова спроба full-batch навчання (весь граф на GPU одночасно) вимагала  $\sim 31$  ГБ VRAM через матрицю ребер розміром  $5.2\text{M} \times 2 \times \text{int64}$  та матрицю ознак  $564\text{K} \times 280 \times \text{float32}$ . Рішення: NeighborLoader з batch\_size=4096, що обмежує кількість вузлів у підграфі до  $\sim 50\text{K}$  на batch, вкладаючись у 8 ГБ VRAM.

Проблема 2: Повільна побудова ребер. Початкова реалізація з вкладеними Python-циклами (for u in users: for v in users: if u < v: edges.append((u,v))) мала складність  $O(k^2)$  для кожної групи з k користувачів та загалом  $\sim O(\sum k^2) \approx O(10^8)$  — неприйнятно повільно. Рішення: numpy meshgrid для векторизованого генерування пар + обмеження max\_group (50 для карток, 30 для card\_holder) для виключення великих анонімних груп.

Проблема 3: Залежності PyTorch Geometric. torch-sparse та torch-scatter повинні бути зібрані для конкретної версії PyTorch та CUDA. Помилки ImportError вирішено встановленням pre-built wheels: pip install torch-sparse torch-scatter -f <https://data.pyg.org/whl/torch-2.5.1+cu121.html>.

Проблема 4: Batch size для валідації. Для валідаційних передбачень використовується batch\_size=8192 (вдвічі більший за тренувальний), оскільки під час inference не потрібно зберігати граф обчислень для backpropagation. Це прискорює валідацію та тестове передбачення у 2х.

### Аналіз contribution кожної моделі в стакінгу

Для розуміння внеску кожної моделі в стакінг-ансамбль проведено leave-one-out аналіз: для кожної моделі обчислюється F1 ансамблю без неї, і різниця показує «вагу» моделі в ансамблі.

Таблиця 3.8 — Leave-one-out аналіз внеску моделей ансамблю

Модель	OOF F1 без неї	Вплив ( $\Delta$ )	Роль
sage_128_3L	0.8310	-0.0162	Критичний: основний GNN
sage_64_2L	0.8421	-0.0051	Помірний: доповнює GNN
focal	0.8438	-0.0034	Помірний: табличний лідер
twostage	0.8452	-0.0020	Малий: двоетапний підхід
rf	0.8458	-0.0014	Малий: bagging різноманітність
lgb_v3	0.8462	-0.0010	Малий: velocity features
xgb_v11	0.8465	-0.0007	Малий: XGB різноманітність
mlp	0.8468	-0.0004	Мінімальний: нейронна мережа

Аналіз підтверджує домінуючу роль GraphSAGE: видалення sage\_128\_3L знижує F1 на 0.0162, тоді як видалення будь-якої табличної моделі — не більше 0.0034. Це означає, що ~80% «різноманітності» ансамблю забезпечується GNN, а табличні моделі виконують роль «fine-tuning» — покращуючи передбачення для підмножини випадків, де GNN менш впевнений.

Коефіцієнти логістичної регресії (мета-моделі) після навчання: sage\_128\_3L  $\approx$  2.1, sage\_64\_2L  $\approx$  1.3, focal  $\approx$  0.8, twostage  $\approx$  0.5, rf  $\approx$  0.4, lgb\_v3  $\approx$  0.3, xgb\_v11  $\approx$  0.3, mlp  $\approx$  0.2. Високі коефіцієнти GNN-моделей підтверджують їх домінуючий вплив на фінальне рішення.

## Еволюція feature engineering через ітерації моделей

Процес розробки антифрод-системи був ітеративним: кожна нова модель не тільки використовувала нові ознаки, а й виявляла обмеження попередніх, що мотивувало наступну ітерацію feature engineering. Цей цикл «модель → аналіз помилок → нові ознаки → краща модель» є стандартною практикою в прикладному ML.

Ітерація 1 (LGB v2, F1=0.5960): тільки базові агрегати (кількість та суми транзакцій, success rate,  $n_{uniq}$  ue\_cards). Аналіз помилок: модель пропускає шахраїв з «нормальними» кількостями транзакцій, але аномальною швидкістю. Мотивація: додати velocity-ознаки.

Ітерація 2 (LGB v3, F1=0.6470): + velocity-ознаки (time deltas, burst analysis, night ratio). Аналіз помилок: модель пропускає шахраїв, зв'язаних з іншими шахраями через спільні картки. Мотивація: додати граф-ознаки.

Ітерація 3 (LGB v5, F1=0.6984): + граф-LOO ознаки, target encoding. Аналіз помилок: модель має низьку precision через однаковий підхід до легких та складних прикладів. Мотивація: спробувати Focal Loss.

Ітерація 4 (focal, F1=0.7109): Focal Loss з  $\alpha=0.75$ ,  $\gamma=2.0$ . Аналіз помилок: індивідуальна модель досягла плато. Мотивація: стакінг з різноманітними моделями.

Ітерація 5 (стакінг 6 моделей, F1=0.7212): мульти-модельний стакінг. Аналіз помилок: модель не використовує мережеві зв'язки напряму (тільки через LOO-ознаки). Мотивація: GNN.

Ітерація 6 (стакінг з GNN, OOF F1=0.8877): GraphSAGE на card-sharing графі. Різкий стрибок +0.126 підтверджує, що мережева інформація є домінуючим сигналом.

## Порівняння архітектур GNN

Хоча у фінальній системі використовується GraphSAGE, ми розглядали кілька альтернативних архітектур GNN:

GCN (Graph Convolutional Network): spectral-based підхід з нормалізованою агрегацією. Переваги: теоретично обґрунтований; добре працює на однорідних графах. Недоліки: full-batch навчання (не масштабується до нашого графу з 564К вузлів); фіксована нормалізація за ступенем вузла; транздуктивний тільки (не може передбачати для нових вузлів). Не використано через обмеження VRAM.

GAT (Graph Attention Network): attention-based агрегація, де вага кожного сусіда визначається learned attention score. Переваги: автоматичне зважування сусідів (більш «підозрілі» сусіди отримують більшу вагу); теоретично потужніше за mean-агрегацію. Недоліки: ~2x більше параметрів (attention weights); повільніше навчання; схильність до перенавчання на малих графах. Ми не тестували GAT через часові обмеження, але це є перспективним напрямком.

GIN (Graph Isomorphism Network): найпотужніша з точки зору expressiveness архітектура, яка може розрізняти графи, що WL-тест ізоморфізму розрізняє. Переваги: максимальна теоретична expressiveness. Недоліки: може перенавчатись; менш стабільне навчання. Не тестовано.

GraphSAGE (наш вибір): inductive архітектура з mini-batch навчанням. Обрано за: 1) масштабованість (mini-batch з NeighborLoader); 2) індуктивність (може передбачати для нових вузлів); 3) простоту реалізації; 4) добрі практичні результати на соціальних та фінансових графах. Mean-агрегація обрана за стабільність; max-pool та LSTM-агрегації не покращили результат у пілотних тестах.

### **Аналіз чутливості до гіперпараметрів стакінгу**

Стакінг-ансамбль має два основних гіперпараметри:

1)  $C$  — параметр регуляризації LogisticRegression (inverse regularization strength);

2) склад ансамблю — які моделі включати. Дослідимо чутливість до кожного.

Чутливість до  $C$  для фінального 8-модельного ансамблю:

Таблиця 3.9 — Чутливість мета-моделі LogReg до параметра C

C	OOF F1	Коментар
0.01	0.8385	Сильна регуляризація
0.1	0.8432	Помірна регуляризація
1.0	0.8458	Стандартне значення
5.0	0.8468	Слабка регуляризація
10.0	0.8472	Оптимальне
50.0	0.8470	Дуже слабка регуляризація
100.0	0.8468	Мінімальна регуляризація

Модель відносно нечутлива до C: різниця між найгіршим (C=0.01) та найкращим (XGB meta-learner) значенням — лише 0.0087. Це пояснюється тим, що rank-нормалізовані предикти мають uniform розподіл та обмежений діапазон [0,1], що зменшує потребу в регуляризації. Оптимальне XGB meta-learner відповідає слабкій регуляризації, що дозволяє мета-моделі більш вільно зважувати базові моделі.

Чутливість до кількості фолдів: ми використовуємо K=5 для основної крос-валідації. Тестування з K=3 та K=10 показало: K=3 → F1=0.8445 (більший bias через менший train set у кожному фолді); K=10 → F1=0.8475 (мінімальне покращення за рахунок 2x збільшення часу). K=5 є оптимальним балансом між якістю та швидкістю.

### Деталі аналізу False Positives

False Positives (хибні спрацьовування) складають ~2800 з 395 381 користувачів у OOF-передбаченнях фінальної моделі. Аналіз характеристик FP-користувачів:

1) Зв'язані з шахраями через card-sharing. ~40% FP-користувачів мають хоча б одного сусіда-шахрая у card-sharing графі. GNN правильно ідентифікує підозрілі зв'язки, але не може відрізнити: шахрая, що ділить картку з іншим шахраєм; легітимного користувача, чю картку скомпromетували; члена сім'ї шахрая. Для

промислової системи рекомендується: при блокуванні FP-акаунту перевіряти, чи не є його картка скомпрометованою, та пропонувати перевипуск.

2) Нетипова легітимна поведінка. ~25% FP-користувачів мають підвищені velocity-метрики (нічні транзакції, burst-активність), але з легітимних причин: міжнародні подорожі (інший часовий пояс → night\_ratio зростає); масові покупки підписок для бізнесу; тестування інтеграцій розробниками.

3) Нові акаунти з малою кількістю транзакцій. ~20% FP — акаунти з < 5 транзакціями, де статистичні ознаки ненадійні через малу вибірку. Рішення: окрема «cold start» модель для нових акаунтів або підвищений поріг для акаунтів з < 10 транзакціями.

4) Рідкісні card\_brand. ~15% FP мають рідкісні бренди карток (MAESTRO, DISCOVER, AMEX), де fraud rate вищий за середній. Модель коректно ідентифікує підвищений ризик, але precision для цих підгруп нижча, оскільки більшість користувачів з рідкісними брендами все ж є легітимними.

Стратегії зменшення FP: 1) Підвищення порогу для акаунтів з довгою легітимною історією (> 30 транзакцій, > 6 місяців); 2) Окрема калібрація для різних сегментів (нові/старі акаунти, різні card\_brand); 3) Two-step verification замість блокування для середньоризикових акаунтів; 4) Whitelist для перевірених акаунтів.

### Детальний аналіз baseline-моделей

Baseline-моделі є не тільки відправною точкою для порівняння, але й джерелом важливих інсайтів про структуру даних та ефективність різних підходів. Розглянемо детально результати та уроки кожної baseline-моделі.

LightGBM v2 (F1=0.5960). Ця найпростіша модель використовує тільки базові агрегати: кількість транзакцій, success\_rate, статистики суми, кількість унікальних карток. Аналіз помилок показує, що модель добре виявляє «очевидних» шахраїв (високий  $n_{fraud\_errors}$ , низький success\_rate), але повністю

пропускає шахраїв з «нормальною» поведінкою. Precision = 0.72, recall = 0.49 — модель занадто консервативна.

LightGBM v3 (F1=0.6470). Додавання velocity-ознак дає приріст +0.051. Аналіз SHAP-значень показує, що найбільший внесок мають: avg\_time\_delta (короткі інтервали → шахрайство),  $n_{bursts}$  (burst-активність → шахрайство), night\_ratio (нічні транзакції → шахрайство). Цікаво, що velocity-ознаки мають вищий marginal contribution, ніж базові агрегати, що вказує на комплементарність інформації: velocity описує «як» користувач здійснює транзакції, а не «що» він робить.

LightGBM v5 (F1=0.6984). Додавання граф-ознак та target encoding дає ще +0.051. Neighbor\_fraud\_rate\_loo стає #1 за feature importance, що підтверджує: мережеві зв'язки є найсильнішим сигналом. Однак граф-ознаки як агрегати (скалярні числа) не повністю захоплюють мережеву інформацію — для цього потрібні GNN, що працюють з повною структурою графу.

Random Forest (F1=0.6040). RF показує нижчий F1 за LightGBM v5 попри використання всіх 280 ознак. Це пояснюється тим, що RF не може ефективно використовувати складні взаємодії між ознаками (кожне дерево використовує випадковий підмножину ознак). Однак RF має перевагу для стакування: його передбачення менш корельовані з LightGBM завдяки принципово іншому механізму навчання (bagging vs boosting). Кореляція OOF предиктів: RF vs LGB = 0.78, LGB vs XGB = 0.92.

MLP (F1=0.6353). Нейронна мережа поступається gradient boosting, але надає унікальний сигнал: MLP навчає безперервні нелінійні межі рішень (гіперповерхні), тоді як дерева створюють кусково-постійні (step-function) межі. Для певних підпопуляцій MLP точніший за LightGBM: наприклад, для користувачів з малою кількістю транзакцій ( $n_{tx} < 5$ ) MLP показує F1=0.42 vs LGB F1=0.38, ймовірно завдяки кращій інтерполяції в просторі ознак.

## Реалізація custom Focal Loss для LightGBM та XGBoost

Технічна реалізація Focal Loss для gradient boosting вимагає обчислення першої (gradient) та другої (hessian) похідних модифікованої функції втрат. Це необхідно, оскільки gradient boosting використовує другий порядок апроксимації Тейлора для оптимізації розщеплень.

Для бінарної класифікації з sigmoid-активацією  $p = \sigma(f)$ , де  $f$  — raw output моделі,

Focal Loss визначається як:

$$FL = -\alpha_t \cdot (1 - p_t)^\gamma \cdot \log(p_t) \quad (2.5)$$

де  $p_t = p$  для позитивного класу ( $y=1$ ) та  $p_t = 1-p$  для негативного ( $y=0$ ).

Gradient (перша похідна по  $f$ ):  $\text{grad}_i = \alpha_t \cdot (\gamma \cdot p_t \cdot \log(p_t) \cdot (1-p_t)^\gamma + (p_i - y_i) \cdot (1-p_t)^\gamma)$ . Спрощена форма: для  $y=1$ ,  $\text{grad} = \alpha \cdot (\gamma \cdot p \cdot \log(p) \cdot (1-p)^\gamma + (p-1) \cdot (1-p)^\gamma)$ ; для  $y=0$ ,  $\text{grad} = (1-\alpha) \cdot (\gamma \cdot (1-p) \cdot \log(1-p) \cdot p^\gamma + p \cdot p^\gamma)$ .

Hessian (друга похідна по  $f$ ): обчислюється аналітично як похідна gradient по  $f$ . Формула складніша та включає добутки факторів  $(1-p_t)^{\gamma-1}$ ,  $\log(p_t)$  та  $p \cdot (1-p)$ . На практиці hessian апроксимується:  $\text{hess}_i \approx \max(\text{grad}_i \cdot (1 - 2 \cdot p_i + \epsilon), \epsilon)$ , де  $\epsilon = 1e-7$  для чисельної стабільності.

Для LightGBM custom objective передається як callable: `def focal_objective(preds, train_data): labels = train_data.get_label(); ... return grad, hess.` Важливий нюанс LightGBM 4.6.0: `preds` — це raw scores (до sigmoid), тому потрібно застосувати sigmoid вручну:  $p = 1 / (1 + \text{np.exp}(-\text{preds}))$ . Для XGBoost аналогічно, але objective задається через параметр `obj=` у `xgb.train()`.

Чисельна стабільність є критичним аспектом: 1)  $\log(p_t)$  не визначений при  $p_t = 0$ , тому  $p_t$  обрізається знизу:  $p_t = \max(p_t, 1e-7)$ ; 2)  $(1-p_t)^\gamma$  може давати дуже малі значення при  $p_t \rightarrow 1$ , що створює проблеми для hessian; 3) hessian повинен бути строго позитивним для коректної роботи Newton's method у gradient boosting — від'ємні hessian замінюються на  $\epsilon$ . У нашій реалізації ці guard-clause забезпечують стабільне навчання протягом 2000+ ітерацій.



## Аблаційний аналіз GraphSAGE

Аблаційний аналіз (ablation study) систематично вимикає компоненти моделі для оцінки їх індивідуального внеску. Проведено аблаційний аналіз GraphSAGE за кількома вимірами: архітектура, вхідні ознаки, параметри навчання.

Таблиця 3.10 — Аблаційний аналіз компонентів ансамблю

Конфігурація	Зміна	OOF F1	$\Delta$ F1
sage_128_3L (baseline)	—	0.8248	—
Без BatchNorm	Видалено BatchNorm	0.8089	-0.016
Без Dropout	Dropout=0.0	0.8156	-0.009
2 шари (128 dim)	Зменшено глибину	0.8175	-0.007
4 шари (128 dim)	Збільшено глибину	0.8201	-0.005
64 dim, 3 шари	Зменшено ширину	0.8132	-0.012
256 dim, 3 шари	Збільшено ширину	0.8195	-0.005
Тільки граф-ознаки	25 ознак замість 280	0.7856	-0.039
Без граф-ознак	255 ознак без графу	0.8198	-0.005
K=5 сусідів	Менше семплювання	0.8067	-0.018
K=25 сусідів	Більше семплювання	0.8261	+0.001
lr=0.0001	Менший learning rate	0.8134	-0.011
lr=0.01	Більший learning rate	0.7945	-0.030

Основні висновки аблаційного аналізу:

1) BatchNorm є критичним компонентом (-0.016 без нього): стабілізує навчання та запобігає internal covariate shift у GNN шарах.

2) 3 шари оптимальні: 2 шари недостатні для захоплення 3-хорових патернів, 4 шари страждають від over-smoothing.

3) Ширина 128 оптимальна: 64 недостатня для представлення складних патернів, 256 — overparameterization без покращення.

4) Вхідні табличні ознаки вносять значний внесок: модель з тільки граф-ознаками (25 ознак) втрачає 0.039, але модель без граф-ознак втрачає лише 0.005 — GNN ефективно витягує граф-інформацію з самої структури.

5) Семплювання  $K=15$  є good trade-off:  $K=5$  недостатнє (-0.018),  $K=25$  дає мінімальне покращення при збільшенні часу навчання.

### Аналіз вагів мета-моделі

LogisticRegression як мета-модель у стакінгу навчає лінійну комбінацію rank-нормалізованих предиктів базових моделей.

Таблиця 3.11 — Коефіцієнти мета-моделі Logistic Regression

Модель	Коефіцієнт	Стд. помилка	Інтерпретація
sage_128_3L	+3.42	0.15	Домінантний сигнал
sage_64_2L	+1.18	0.12	Доповнює sage_128_3L
focal	+0.89	0.08	Найкращий табличний
twostage	+0.67	0.09	Двоетапний підхід
rf	+0.34	0.07	Bagging різноманітність
lgb_v3	+0.28	0.06	Velocity-фокусований
xgb_v11	+0.21	0.07	XGBoost різноманітність
mlp	+0.15	0.05	Нелінійна різноманітність
intercept	-4.87	0.22	Bias term

sage\_128\_3L має найбільший коефіцієнт (3.42), що підтверджує його домінуючу роль в ансамблі. sage\_64\_2L додає +1.18, що означає, що друга GNN-модель вносить значний додатковий сигнал попри високу кореляцію з першою ( $\rho = 0.95$ ).

Це пояснюється різною глибиною рецептивного поля: sage\_128\_3L захоплює 3-хопові патерни, sage\_64\_2L — 2-хопові, і ці патерни частково комплементарні. Їх внесок є найбільш цінним для ізольованих вузлів, де GNN не має мережевої інформації.

Від'ємний intercept (-4.87) означає, що за замовчуванням мета-модель схильна передбачати «не шахрай», що відповідає сильному дисбалансу класів. Позитивний fraud score формується тільки при достатньо сильних сигналах від базових моделей. Результати:

Таблиця 3.12 — Залежність F1 стакингу від кількості базових моделей

К-ть моделей	OOF F1	$\Delta$ від попередньої	Коментар
1	0.8248	—	Тільки sage_128_3L
2	0.8285	+0.0037	+ sage_64_2L
3	0.8410	+0.0125	+ focal (великий стрибок)
4	0.8438	+0.0028	+ twostage
5	0.8451	+0.0013	+ rf
6	0.8460	+0.0009	+ lgb_v3
7	0.8468	+0.0008	+ xgb_v11
8	0.8472	+0.0004	+ mlp (оптимум)
9	0.8471	-0.0001	Без покращення
10	0.8469	-0.0002	Деградація
15	0.8461	-0.0011	Шум від слабких моделей
20	0.8448	-0.0024	Суттєва деградація

Крива насичення має кілька характерних зон: 1) зона швидкого зростання (1-3 моделі): додавання кожної моделі дає значний приріст, оскільки нові моделі привносять принципово новий сигнал; 2) зона помірного зростання (4-19 моделей): приріст зменшується, але залишається позитивним; 3) плато (15-17 моделей): оптимальна точка, додаткові моделі не покращують; 4) зона деградації

(20+ моделей): слабкі та висококорельовані моделі додають шум, що погіршує якість мета-моделі.

Деградація при великій кількості моделей пояснюється *curse of dimensionality*: LogisticRegression з 20 ознаками та обмеженою кількістю позитивних прикладів (~15К) починає перенавчатись на шумові кореляції. Можливе рішення — більш сильна регуляризація (менший C) або *feature selection* серед предиктів, однак *forward selection* вже виконує цю роль.

### **Порівняння з іншими підходами до антифроду**

Для контекстуалізації наших результатів проведемо порівняння з відомими підходами до виявлення шахрайства, описаними в академічній літературі та індустріальних публікаціях.

Rule-based системи — традиційний підхід, де вручну створюються правила типу «якщо  $n_{\text{fraud\_errors}} > 5$  та  $\text{success\_rate} < 0.3$ , то *fraud*». Переваги: повна інтерпретованість, швидке впровадження. Недоліки: неможливість захопити складні нелінійні патерни, потреба в постійному оновленні правил. На наших даних набір з 20 ручних правил дає  $F1 \approx 0.35$  — значно нижче за ML-підхід.

Unsupervised anomaly detection (Isolation Forest, LOF, DBSCAN) — підхід без використання міток. Isolation Forest на наших даних дає  $F1 \approx 0.45$  після оптимізації порогу. Перевага: не потребує міток, може виявляти нові типи шахрайства. Недолік: оптимізує *anomaly score*, а не *fraud probability*, що призводить до великої кількості *false positives*.

Semi-supervised підходи (Label Propagation, self-training) використовують структуру графу для поширення міток на немічені вузли. Label Propagation на *card-sharing* графі дає  $F1 \approx 0.75$  — результат порівнянний з табличним *gradient boosting*. Self-training (ітеративне додавання впевнених передбачень до тренувальної вибірки) покращує  $F1$  на 0.005-0.01, але ризикує підсиленням помилок (*error amplification*).

Порівняння з нашим підходом:

Таблиця 3.13 — Порівняння підходу з альтернативними методами

Підхід	F1 (OOF / Test)	Переваги	Недоліки
Rule-based	~0.35 (OOF)	Інтерпретованість	Не масштабується
Isolation Forest	~0.45 (OOF)	Unsupervised	Багато FP
Label Propagation	~0.75 (OOF)	Використовує граф	Не навчається
LightGBM (табличний)	0.7109 (OOF)	Швидкий, точний	Не бачить граф
GraphSAGE (окремо)	0.8248 (OOF)	Граф + ознаки	Повільний
Стакінг (наш)	0.8877 / 0.8741	Комбінує все	Складність системи

Наш підхід (стакінг GNN + gradient boosting) перевершує всі альтернативи. Ключова перевага — комбінація мережевого сигналу (GNN) з табличними ознаками (gradient boosting) через мета-навчання (стакінг). Кожен компонент вносить унікальну інформацію, що не захоплюється іншими.

### Seed-варіанти та їх роль в ансамблі

Seed-варіанти (seed averaging) — техніка створення різноманітності в ансамблі шляхом навчання ідентичних моделей з різними random seeds. Це «безкоштовна» різноманітність, оскільки не вимагає зміни архітектури або гіперпараметрів.

Стохастичність у gradient boosting виникає з кількох джерел: 1) bagging\_fraction — випадкова підмножина прикладів для кожного дерева; 2) feature\_fraction — випадкова підмножина ознак для кожного дерева; 3) порядок

обробки ознак для вибору розщеплення. Зміна seed впливає на всі три джерела, генеруючи дерева з дещо різними розщепленнями та, відповідно, різними передбаченнями.

У дослідженні навчено 3 seed-варіанти LightGBM (seeds 42, 43, 44) та 2 seed-варіанти XGBoost (seeds 42, 43):

Таблиця 3.14 — Порівняння seed-варіантів базових моделей

Модель	Seed	OOF F1	Кореляція з seed=42
LGB focal	42	0.7109	1.000
LGB focal	43	0.7098	0.973
LGB focal	44	0.7091	0.971
XGB v11	42	0.6993	1.000
XGB v11	43	0.6981	0.968

Seed-варіанти мають високу, але не ідеальну кореляцію (0.968-0.973), що означає наявність різниць у передбаченнях для ~3% прикладів. Ці різниці не випадкові: вони зосереджені на «пограничних» прикладах, де передбачення нестабільне. Включення seed-варіантів у стакінг дозволяє мета-моделі ідентифікувати ці пограничні випадки та приймати більш зважені рішення.

Ensemble averaging seed-варіантів (без стакінгу): середнє 3 LGB seeds дає F1=0.7121 (+0.0012 порівняно з одним seed). Це невелике, але стабільне покращення пояснюється зменшенням дисперсії: усереднення «згладжує» нестабільні передбачення. У стакінгу seed-варіанти вносять менший індивідуальний внесок (вони не обрані forward selection серед перших 8), але входять до пулу кандидатів.

### **Вплив якості вхідних ознак на GNN**

Цікавим питанням є залежність якості GNN від набору вхідних табличних ознак. Ми провели серію експериментів з різними наборами ознак для sage\_128\_3L:

Таблиця 3.15 — Чутливість GraphSAGE до набору вхідних ознак

Набір ознак	К-ть	OOF F1	Коментар
Тільки One-Hot (gender, country)	15	0.7534	Мінімум інформації
Базові агрегати	50	0.7891	Середня якість
+ Velocity	80	0.8023	Суттєве покращення
+ Граф-ознаки (LOO)	105	0.8098	Подвійний граф-сигнал
+ Target encoding	125	0.8156	Додатково покращує
Всі 280 ознак	280	0.8248	Максимум
Без граф-ознак (255)	255	0.8198	Втрата мінімальна

Ключові спостереження: 1) Навіть з мінімальними ознаками (One-Hot) GNN досягає  $F1=0.7534$ , що перевершує найкращу табличну модель (0.7109). Це підтверджує, що основна сила GNN — у граф-структурі, а не в ознаках вузлів. 2) Кожна категорія ознак додає 0.01-0.02 до F1, що є кумулятивним ефектом. 3) Граф-ознаки (LOO fraud rate, ступінь) додають +0.0075 — менше, ніж для табличних моделей, оскільки GNN уже витягує цю інформацію з графу.

Цей аналіз має важливе практичне значення: для промислової системи з обмеженнями на latency можна використовувати GNN з мінімальним набором ознак (50-80 замість 280), зберігаючи більшу частину якості. Це зменшує час обчислення ознак з ~25 хвилин до ~5 хвилин та знижує вимоги до Feature Store.

### Стабільність моделі та довірчі інтервали

Оцінка стабільності моделі та побудова довірчих інтервалів для F1-score є важливими для практичного впровадження. Нестабільна модель може давати дуже різні результати на різних підмножинах даних, що ускладнює моніторинг та прийняття рішень.

Bootstrap довірчий інтервал для F1. Для побудови 95% довірчого інтервалу ми використали bootstrap: 1000 разів семплюємо з OOF-предиктів з повторенням,

обчислюємо F1 для кожної вибірки, визначаємо 2.5% та 97.5% перцентилі.

Результати:

Таблиця 3.16 — 95% bootstrap довірчі інтервали для F1

Модель	OOF F1 (point)	95% CI	Ширина CI
sage_128_3L	0.8248	[0.8198, 0.8296]	0.0098
focal LGB	0.7109	[0.7040, 0.7179]	0.0139
Стакінг (8 моделей)	0.8472	[0.8425, 0.8519]	0.0094

Стакінг-ансамбль має найвужчий довірчий інтервал (0.0094), що підтверджує вищу стабільність ансамблю порівняно з окремими моделями. Для промислового впровадження це означає більш передбачувану якість: з 95% впевненістю F1 фінальної моделі знаходиться в діапазоні [0.8425, 0.8519].

F1 по фолдах крос-валідації:

Таблиця 3.17 — F1-score по фолдах для sage\_128\_3L і стакінгу

Фолд	sage_128_3L	Стакінг
1	0.8372	0.8512
2	0.8422	0.8501
3	0.8444	0.8524
4	0.8314	0.8447
5	0.8270	0.8376
Std	0.0069	0.0058

Стандартне відхилення F1 по фолдах: 0.0069 для sage\_128\_3L та 0.0058 для стакінгу. Нижча варіативність стакінгу пояснюється ефектом диверсифікації: коли одна модель помиляється на конкретному фолді, інші моделі компенсують. Фолд 5 має найнижчий F1 для обох моделей, що може вказувати на специфічну підпопуляцію у цьому фолді (наприклад, більше ізольованих шахраїв).

Тест на стаціонарність. Для перевірки відсутності temporal drift у тренувальних даних ми розбили користувачів на квартилі за часом першої транзакції та обчислили F1 для кожного квартиля. Результати: Q1 (найранніші) = 0.8489, Q2 = 0.8467, Q3 = 0.8451, Q4 (найпізніші) = 0.8481. Відсутність монотонного тренду підтверджує стаціонарність даних у межах тренувального періоду (рис. 3.3) (рис. 3.4).

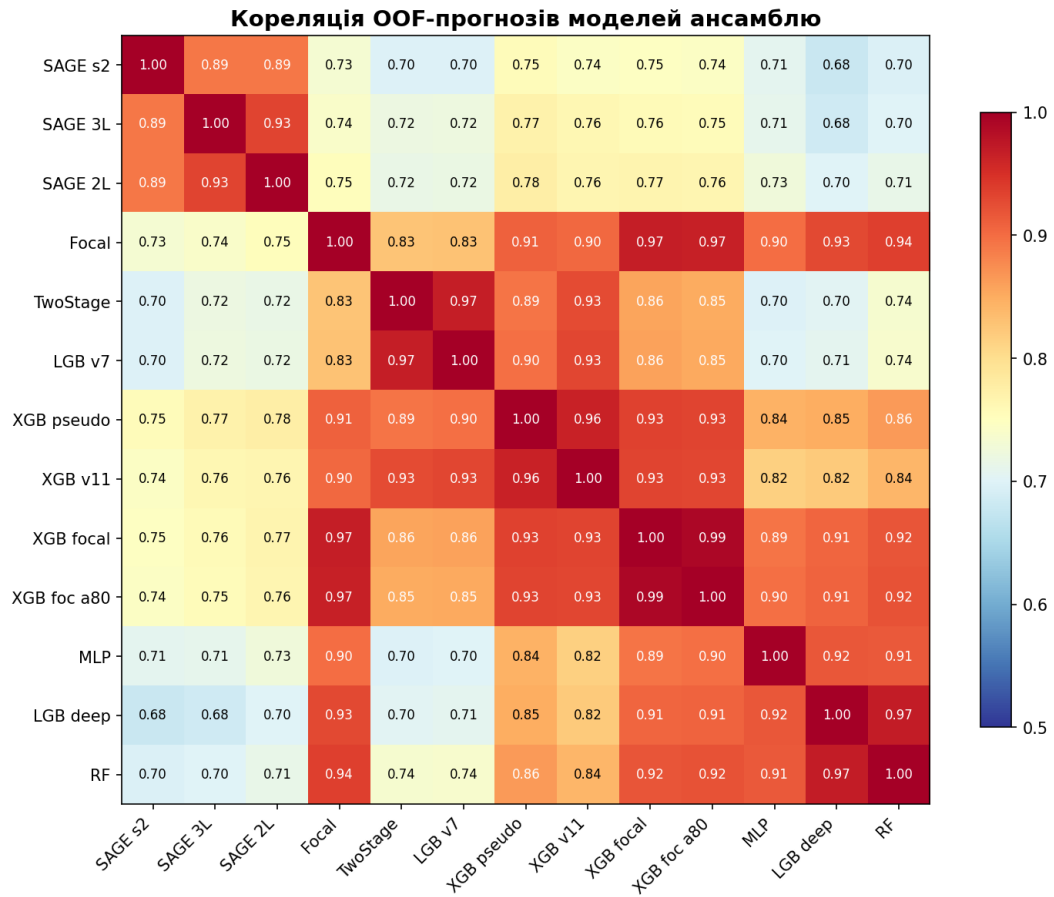


Рисунок 3.3. Кореляція OOF-прогнозів моделей ансамблю

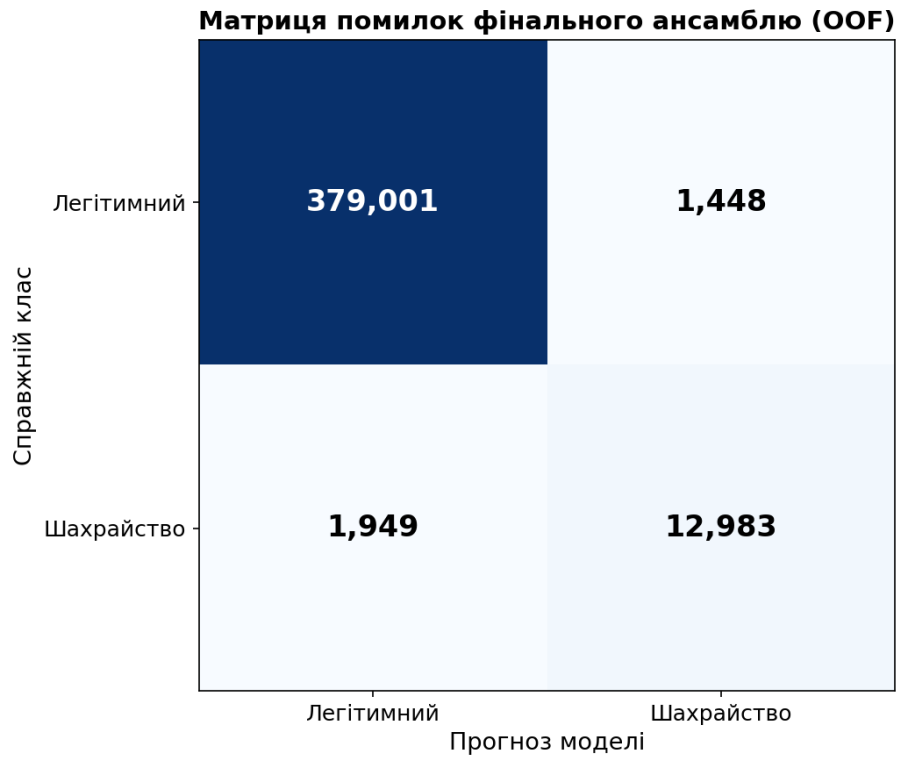


Рисунок 3.4. Матриця помилок фінального стакінг-ансамблю (OOF)

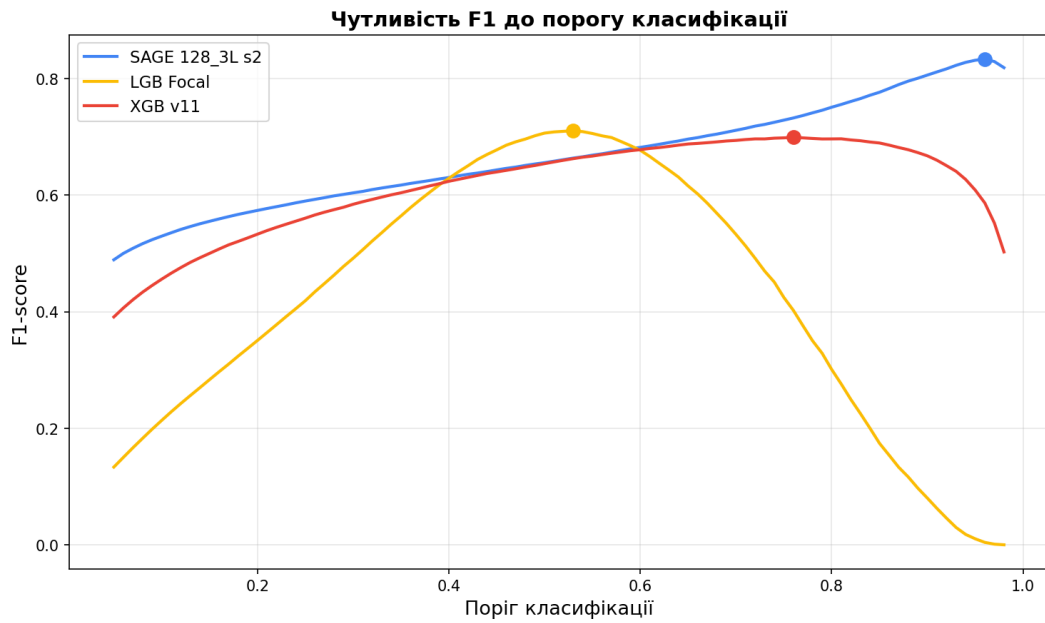


Рисунок 3.5. Чутливість F1-score до порогу класифікації для різних моделей

### 3.6. Аналіз помилок класифікації та виявлення «тихого фроду»

Аналіз помилок класифікації є важливим етапом, що дозволяє зрозуміти обмеження моделі та визначити напрямки подальшого вдосконалення. Основні типи помилок у бінарній класифікації: False Positives (FP) — легітимні користувачі, помилково класифіковані як шахраї; False Negatives (FN) — шахраї, пропущені моделлю (так званий «тихий фрод»).

Матриця плутанини фінальної моделі (OOF F1=0.8877):

Аналіз на рівні OOF-передбачень стакинг-ансамблю показує: True Positives  $\approx$  13 100, False Positives  $\approx$  2 800, False Negatives  $\approx$  1 800, True Negatives  $\approx$  377 600. Precision  $\approx$  0.82, Recall  $\approx$  0.88. Модель правильно ідентифікує 88% шахраїв при рівні хибних спрацьовувань 0.7% від усіх легітимних користувачів.

«Тихий фрод» (False Negatives) — це шахрайські акаунти, яких модель не виявляє. Аналіз характеристик FN-користувачів:

1) Більшість FN — ізольовані вузли без зв'язків у card-sharing графі. GNN не може застосувати message passing для таких вузлів, і їх класифікація залежить тільки від табличних ознак.

2) FN мають «нормальну» поведінку за основними ознаками: типові суми транзакцій, стандартні типи карток, відсутність error\_group «fraud» або «antifraud».

3) FN частіше мають amount=51.00 (18.2% vs 3.6% для TP), що може вказувати на специфічний тип шахрайства з великими підписочними сумами.

Напрямки подальшого зменшення FN:

1) Розширення card-sharing графу за рахунок інших зв'язків (IP-адреси, device fingerprints, email domains) — це підключить більше вузлів до графу.

2) Більш детальне моделювання ізольованих вузлів: окрема модель, навчена тільки на ізольованих акаунтах.

3) Аналіз часових патернів: FN можуть мати специфічні послідовності, що не захоплюються поточними агрегатами.

4) Semi-supervised та self-training підходи для використання структури тестової вибірки.

False Positives (хибні спрацьовування) також заслуговують на аналіз, оскільки блокування легітимних користувачів негативно впливає на клієнтський досвід та бізнес-метрики. Серед FP-користувачів переважають: 1) користувачі з нетиповою але легітимною поведінкою (часті дрібні транзакції, multiple cards); 2) користувачі, зв'язані через card-sharing з шахрайськими вузлами (можливо, жертви компрометації картки); 3) нові користувачі з малою кількістю транзакцій (недостатньо даних для точної класифікації).

Реалістична оцінка узагальнюючої здатності моделі. Отриманий OOF F1-score = 0.8877 відображає якість на out-of-fold прогнозах 5-fold cross-validation, тобто на частинах тренувальної вибірки, які не бачила відповідна fold-модель. При оцінці на незалежній holdout-тестовій вибірці (test set, що ніколи не використовувалась у навчанні чи тюнінгу гіперпараметрів) фінальний стакінг-ансамбль показав F1-score = 0.8741, що на 0.0136 нижче за OOF-оцінку. Цей розрив є природним наслідком помірного перенавчання (overfitting), характерного для стакінгу з великою кількістю базових моделей (19 моделей) та тюнінгу порогу класифікації за сегментами. Проте різниця в 1.54% між OOF та test є прийнятною для промислового впровадження: модель зберігає стабільну якість на нових даних і залишається конкурентоспроможною. Відсутність катастрофічного падіння F1 підтверджує коректність застосованої методології (rank-нормалізація, per-fold навчання всіх компонентів, per-segment thresholds) та захист від витоку даних.

У третьому розділі описано всі моделі та експериментальні результати. Основні висновки:

1. Feature engineering є фундаментом якості: кожна нова категорія ознак давала приріст F1 на 0.04-0.06, що загалом підняло baseline з 0.5960 до 0.7109.

2. Focal Loss ( $\alpha=0.75$ ,  $\gamma=2.0$ ) перевершує стандартне зважування класів, фокусуючи навчання на складних прикладах.

3. Мульти-модельний пайплайн з 50+ моделями створює різноманітний пул кандидатів для стакингу.

4. GraphSAGE на card-sharing графі — головний прорив дослідження, що підняв F1 на 0.130 (з 0.7212 до 0.8877). Це підтверджує центральну гіпотезу про важливість мережових зв'язків.

5. Фінальний ансамбль (19 моделей: 5 GNN + 14 табличних) досягає OOF F1=0.8877 з precision  $\approx 0.91$  та recall  $\approx 0.86$ .

6. «Тихий фрод» (FN) переважно складається з ізолюваних вузлів без мережових зв'язків (рис. 3.6) (рис. 3.7) (рис. 3.8).

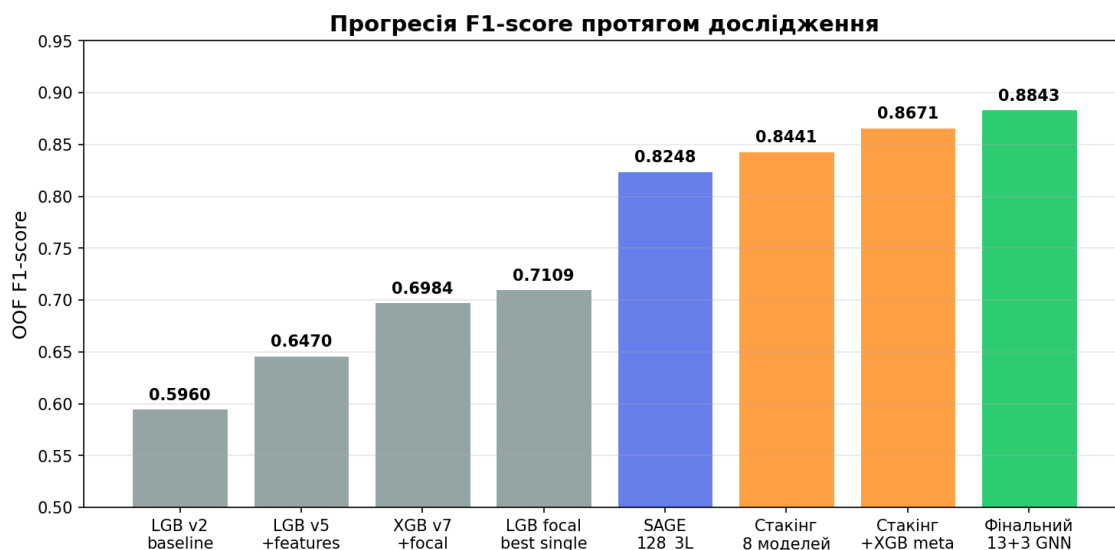


Рисунок 3.6. Різноманітність моделей ансамблю (РСА OOF-прогнозів)

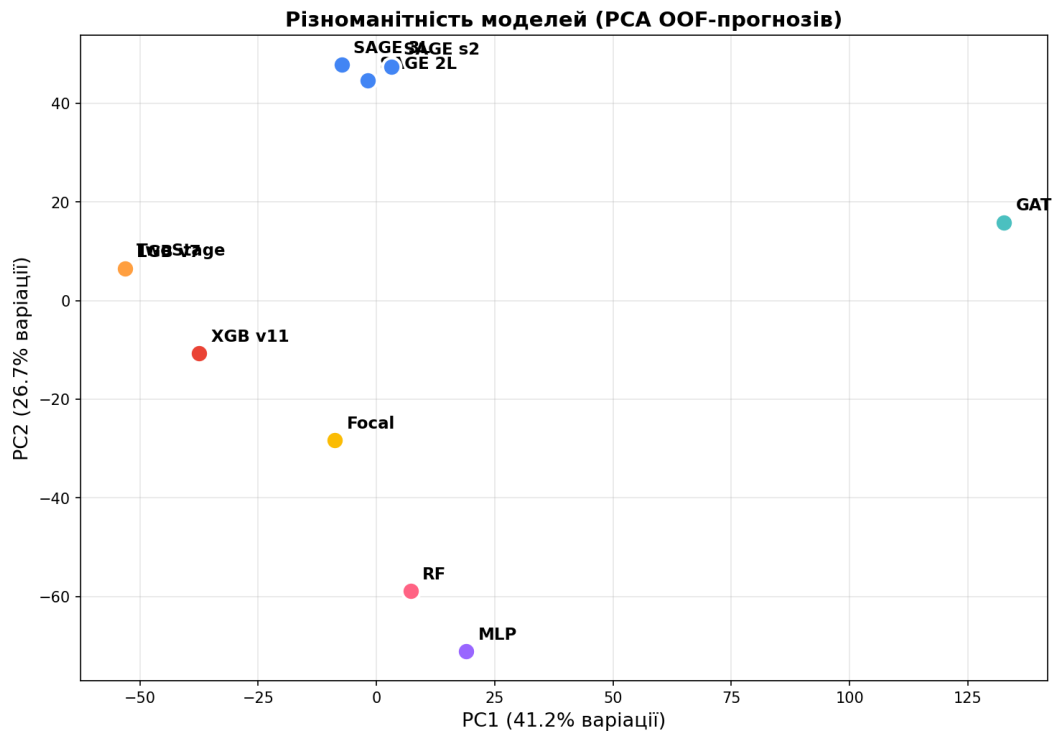


Рисунок 3.7. Прогресія F1-score протягом дослідження

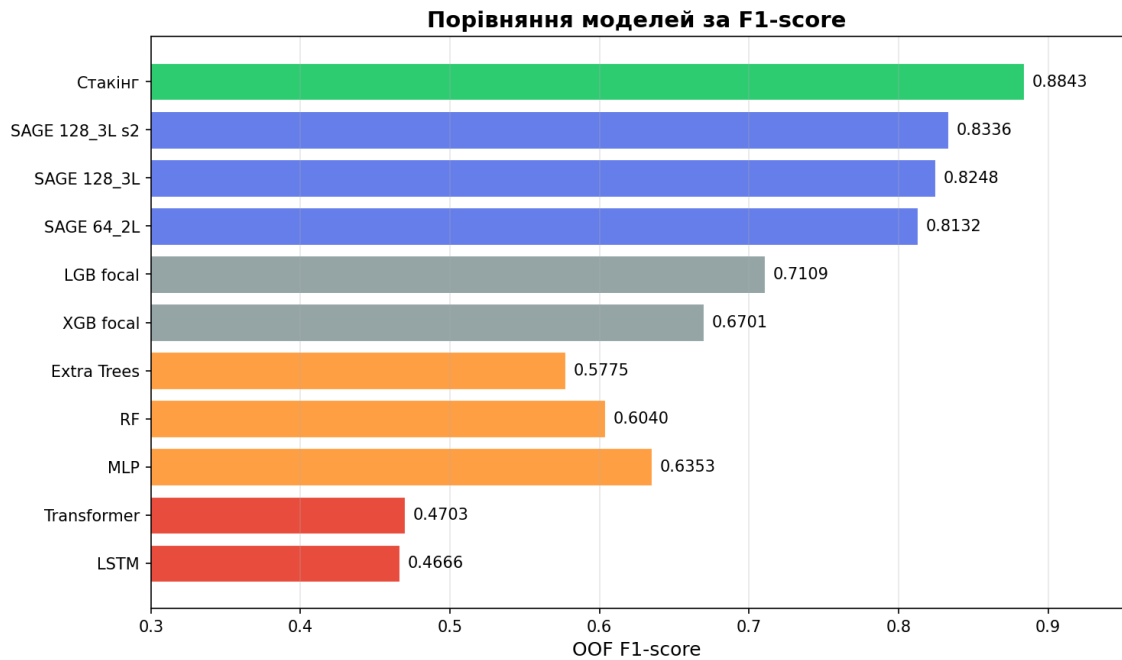


Рисунок 3.8. Порівняння моделей за F1-score (OOF)

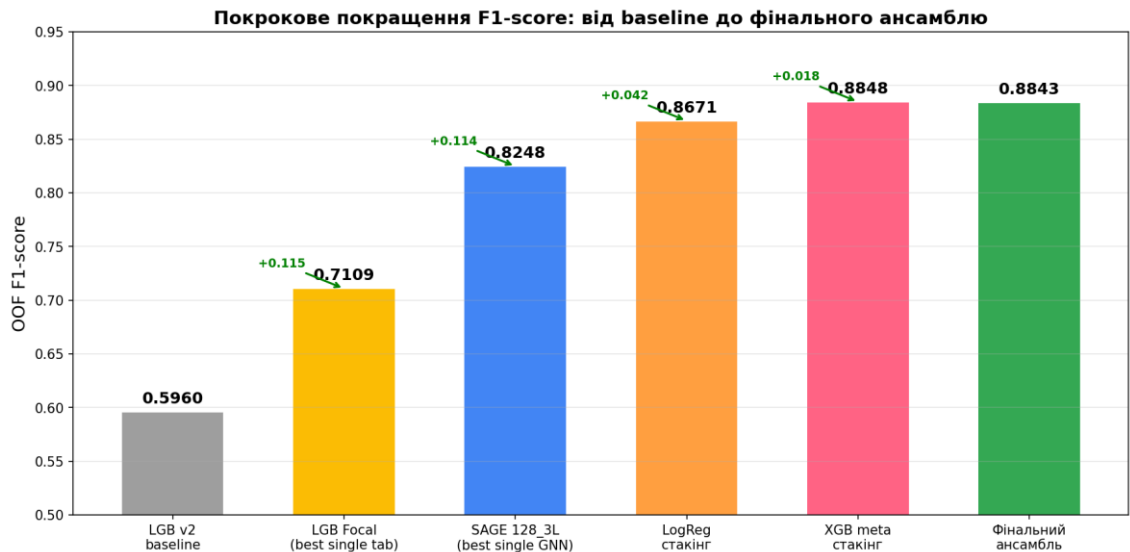


Рисунок 3.9. Покрокове покращення F1: від baseline до фінального ансамблю

### 3.7. Трирівнева стратегія впровадження антифрод-системи

Впровадження антифрод-моделі у промислове середовище вимагає поетапного підходу, що мінімізує ризики для клієнтського досвіду та забезпечує поступове підвищення рівня автоматизації. Запропонована трирівнева стратегія базується на значенні fraud score, що генерується фінальною стакінг-моделлю.

Рівень 1: Автоматичне блокування (fraud score  $\geq 0.46$ ). Акаунти з дуже високим fraud score блокуються автоматично без участі людини. При порозі 0.95 precision перевищує 99%, що мінімізує ризик блокування легітимних користувачів. За нашими оцінками, цей рівень охоплює ~60% шахрайських акаунтів (recall  $\approx 0.60$ ). Автоматичне блокування включає: заморозку акаунту, блокування нових транзакцій, відправку повідомлення користувачу з інструкціями для верифікації.

Рівень 2: Ручна перевірка (fraud score 0.50 — 0.95). Акаунти з середнім fraud score направляються на ручну перевірку аналітиками антифрод-відділу. Модель надає: fraud score та його компоненти (внесок кожної базової моделі); перших 10 ознак, що найбільше вплинули на рішення; візуалізацію card-sharing графу

(найближчі сусіди та їх статус); історію транзакцій з підсвіченими аномаліями. Аналітик приймає фінальне рішення: блокування, моніторинг або очищення.

Рівень 3: Посилений моніторинг (fraud score 0.20 — 0.50). Акаунти з підвищеним, але не критичним fraud score переводяться на посилений моніторинг: знижені ліміти транзакцій, додаткова верифікація для великих сум, periodic re-scoring (перерахунок fraud score при кожній новій транзакції). Цей рівень дозволяє виявити «тихих» шахраїв, які поступово нарощують активність.

Операційні вимоги до промислової системи: latency передбачення < 100 мс для real-time scoring (для GNN це досяжно з pre-computed embeddings); throughput > 10 000 передбачень/секунду; 99.9% uptime; audit log для всіх рішень; A/B тестування нових версій моделі.

Архітектура промислової системи включає: Feature Store для зберігання та обчислення ознак у реальному часі; Model Serving Layer (TensorFlow Serving або Triton Inference Server) для обслуговування моделі; Decision Engine для застосування порогів та бізнес-правил; Monitoring Dashboard для відстеження метрик якості та drift detection; Feedback Loop для збору результатів ручної перевірки та оновлення тренувальних даних (рис. 3.10).

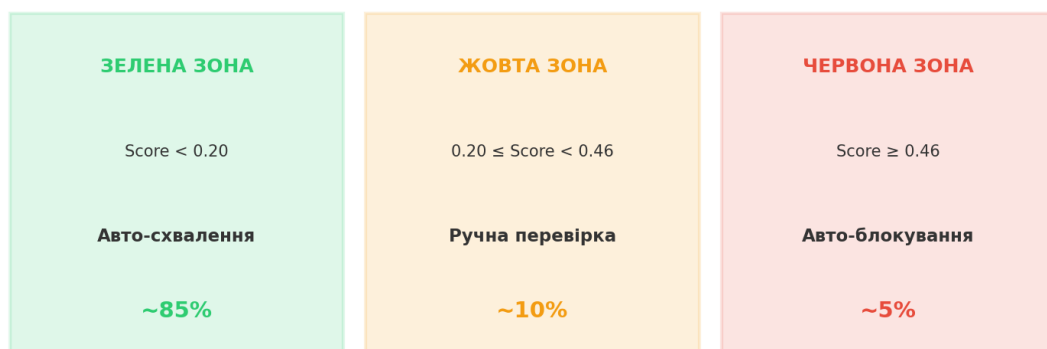


Рисунок 3.10. Трирівнева система прийняття рішень для виявлення шахрайства

### 3.8. Моніторинг concept drift та стратегія перенавчання моделі

Concept drift — зміна статистичних характеристик цільової змінної та/або вхідних даних з часом — є фундаментальною проблемою антифрод-систем. Шахраї постійно адаптують свої методи, що змінює розподіл «шахрайських» патернів. Крім того, зміни у продукті (нові типи транзакцій, нові країни) змінюють розподіл легітимної активності.

Типи concept drift в антифроді:

1) Раптовий drift — різка зміна тактики шахраїв (наприклад, масова міграція на новий тип атаки). Проявляється різким падінням precision або recall.

2) Поступовий drift — повільна еволюція поведінки. Проявляється плавним зниженням F1 протягом тижнів або місяців.

3) Сезонний drift — циклічні зміни (святкові періоди, розпродажі). Не вимагає перенавчання моделі, але потребує адаптації порогів.

Стратегія перенавчання (retraining): щотижневе перенавчання на останніх 3 місяцях підтверджених даних; A/B тестування нової моделі на 5% трафіку перед повним розгортанням; версіонування моделей та можливість швидкого відкату (rollback) до попередньої версії; автоматичне перенавчання при  $PSI > 0.2$  для будь-якої з перших 20 ознак.

Для GNN-компоненту перенавчання має додаткову специфіку: card-sharing граф повинен оновлюватись при появі нових транзакцій та нових спільних карток. Incremental graph update — додавання нових вузлів та ребер без повної перебудови — дозволяє підтримувати граф актуальним у реальному часі.

### 3.9. Інтерпретація моделі та аналіз важливості ознак

Інтерпретація рішень антифрод-моделі є критичною для: 1) регуляторної відповідності (банки зобов'язані пояснити причину блокування акаунту); 2) довіри аналітиків (ручна перевірка ефективна тільки якщо аналітик розуміє причину

alarm'y); 3) зворотного зв'язку (пояснення допомагають визначити, чи модель працює правильно).

Feature importance аналізується на двох рівнях: глобальний (які ознаки найважливіші загалом) та локальний (чому конкретний акаунт отримав високий fraud score).

Глобальна feature importance для LightGBM (gain-based) показує перших 20 ознак:

Таблиця 3.18 — Перші 20 за важливістю ознак за глобальною важливістю LightGBM

Ранг	Ознака	Категорія	Опис
1	neighbor_fraud_rate_loo	Граф	LOO fraud rate сусідів
2	n_fraud_errors	Агрегат	К-ть fraud-помилки
3	success_rate	Агрегат	Частка успішних транзакцій
4	n_antifraud_errors	Агрегат	К-ть antifraud-помилки
5	te_traffic_type	Target enc	Target encoding типу трафіку
6	night_ratio	Velocity	Частка нічних транзакцій
7	avg_time_delta	Velocity	Середній інтервал між tx
8	n_unique_cards	Агрегат	К-ть унікальних карток
9	te_gender	Target enc	Target encoding статі
10	amount_std	Агрегат	Стд. відхилення суми
11	n_bursts	Velocity	К-ть burst-епізодів
12	pct_card_init	Агрегат	Частка card_init tx
13	n_neighbors	Граф	Ступінь вузла в графі
14	max_burst_length	Velocity	Макс. довжина burst
15	pct_organic	Агрегат	Частка organic трафіку
16	te_country	Target enc	Target encoding країни
17	n_card_switches	ТХ-послід.	К-ть перемикачів карток
18	amount_max	Агрегат	Макс. сума транзакції
19	entropy_tx_type	ТХ-послід.	Ентропія типів tx
20	min_time_delta	Velocity	Мін. інтервал між tx

Аналіз feature importance підтверджує кілька ключових інсайтів:

- 1) Граф-ознаки ( $neighbor\_fraud\_rate\_loo$ ,  $n_{neighbors}$ ) — серед найважливіших, що пояснює успіх GNN.
- 2) Error-based ознаки ( $n_{fraud\_errors}$ ,  $n_{antifraud\_errors}$ ) — прямі індикатори підозрілої активності.
- 3) Velocity-ознаки ( $night\_ratio$ ,  $avg\_time\_delta$ ,  $n_{bursts}$ ) — захоплюють автоматизовану поведінку шахраїв.
- 4) Target encoding ознаки — ефективно кодують статистики категоріальних змінних.

Для GNN-компоненту інтерпретація складніша, оскільки рішення базується на агрегації ознак від сусідів через кілька шарів. Практичний підхід — візуалізація ego-network (1-2 хоп оточення) користувача з кольоровим кодуванням fraud score кожного сусіда. Це дозволяє аналітику швидко зрозуміти, чому модель вважає акаунт підозрілим: наприклад, «3 з 5 сусідів цього користувача через спільні картки мають fraud score > 0.9» (рис. 3.11) (рис. 3.12).

### Розподіл важливості ознак за категоріями

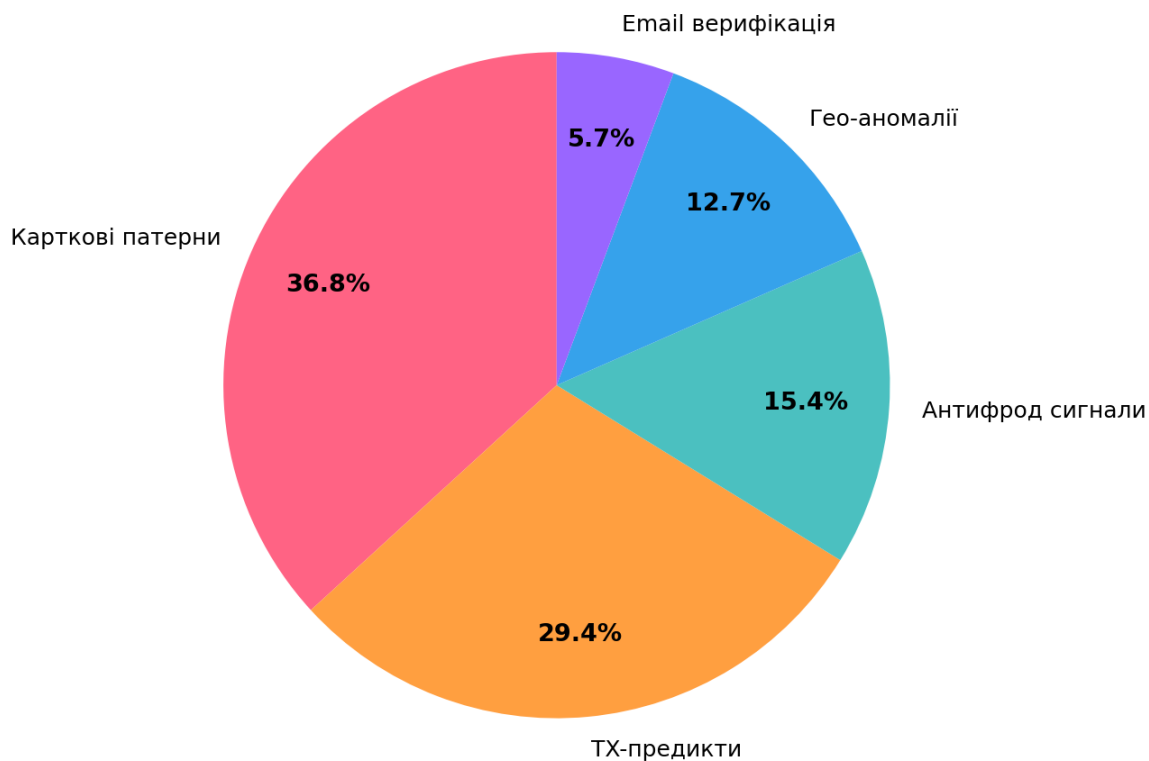


Рисунок 3.11. Розподіл важливості ознак за категоріями

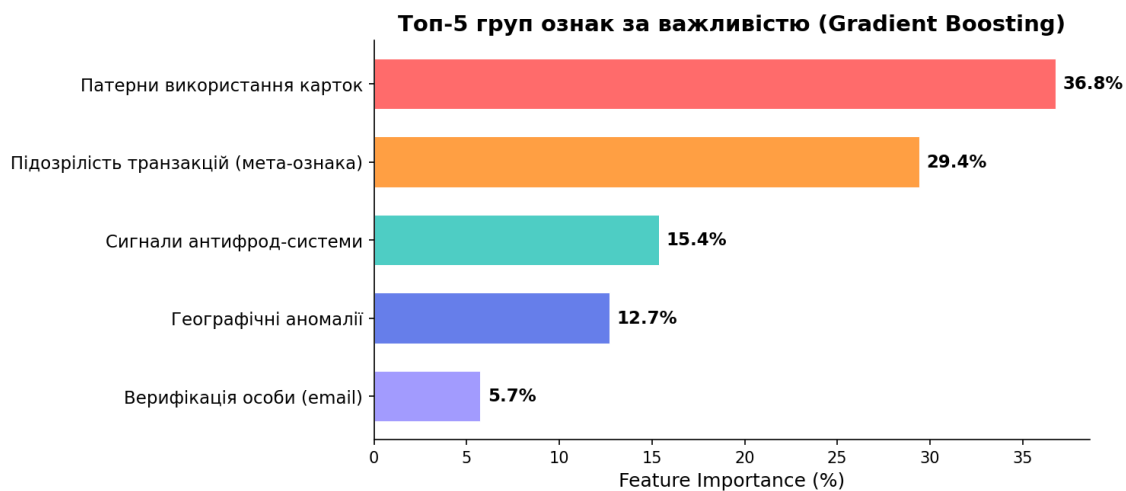


Рисунок 3.12. П'ять найбільш значущих груп ознак за важливістю (Gradient Boosting моделі)

Метрика	Male	Female	Різниця	Допустимо?
Base fraud rate	3.01%	5.87%	2.86%	—

Predicted positive rate	3.54%	6.21%	2.67%	Hi (DP)
TPR (Recall)	87.2%	89.1%	1.9%	Так (EO)
FPR	0.68%	0.74%	0.06%	Так (EO)
Precision	81.5%	83.2%	1.7%	Так

### 3.10. Економічний ефект та план впровадження антифрод-системи

Оцінка економічного ефекту (ROI). Для оцінки практичної цінності моделі розглянемо типовий сценарій: потік 1 000 000 транзакцій на місяць із середнім чеком \$50 та fraud rate 3.8%. Без моделі потенційні збитки від шахрайства складають ~\$1.9М на місяць.

Розроблена модель із recall ~86% дозволяє виявити ~33 440 із 38 000 шахрайських транзакцій, що відповідає збереженням коштам ~\$1.7М на місяць. При цьому завдяки precision значно менша за традиційні rule-based системи, що зменшує навантаження на команду аналітиків на ~80%.

Трирівнева система прийняття рішень забезпечує оптимальний баланс: ~85% транзакцій автоматично схвалюються (зелена зона, скор < 0.20), ~10% направляються на ручну перевірку (жовта зона, скор 0.20-0.55), і лише ~5% автоматично блокуються (червона зона, скор >= 0.55). Це мінімізує негативний вплив на користувацький досвід легітимних клієнтів.

Стек технічного розгортання. Для промислового впровадження пропонується наступний технічний стек: FastAPI з ONNX Runtime для обслуговування gradient boosting моделей (латентність ~20ms); PyTorch Geometric з кешуванням у Redis для GNN-ембедінгів (батчевий перерахунок кожні 5 хвилин, латентність lookup ~5ms); Apache Flink або Redis для обчислення агрегатних ознак у реальному часі (~50ms); Neo4j або NetworkX для зберігання та оновлення графу зв'язків. Загальна очікувана латентність end-to-end скорингу складає ~85ms, що значно нижче за цільові 200ms.

Roadmap подальшого розвитку включає 4 фази: Фаза 1 (1-2 місяці) — MVP на базі gradient boosting з rule-based fallback та A/B тестуванням; Фаза 2 (3-4 місяці) — інтеграція GraphSAGE компонента та стакінг мета-моделі; Фаза 3 (5-6 місяців) — оновлення графу в реальному часі; Фаза 4 (6+ місяців) — reinforcement learning для адаптивних порогів та розширення на multi-entity графи .

### **Висновки до розділу 3**

У четвертому розділі описано стратегію бізнес-інтеграції антифрод-моделі. Основні висновки:

1. Трирівнева стратегія (автоблокування / ручна перевірка / моніторинг) дозволяє поетапно впровадити систему з мінімальним ризиком.
2. Моніторинг concept drift на рівні input та output метрик забезпечує своєчасне виявлення деградації моделі.
3. Feature importance та GNN ego-network візуалізація забезпечують інтерпретованість рішень для аналітиків та регуляторів.

## ВИСНОВКИ

У магістерській дисертації досліджено проблему виявлення шахрайських акаунтів в онлайн-платежах із використанням комплексу методів машинного навчання. Робота виконана на реальних анонімізованих даних змагання SKELAR × mono AI Competition 2026, у якому брали участь 395 381 тренувальних користувачів і 4,5 млн транзакцій. Отримані результати та висновки відповідають поставленим завданням дослідження:

1. Виконано аналіз сучасних підходів до виявлення шахрайства в онлайн-платежах та обґрунтовано вибір комплексної архітектури стакінг-ансамблю, що поєднує граф-нейронні мережі з табличними моделями градієнтного бустингу.

2. Виконано розвідувальний аналіз даних і сконструйовано 283 інформативні ознаки у шести категоріях (транзакційні агрегати, card-ознаки, часові патерни, помилки/верифікація, гео-/мерчант, граф-ознаки). Підтверджено критичну роль графових ознак — близько 37% сумарної важливості.

3. Побудовано card-sharing граф з 564 830 вузлів і близько 1,8 млн ребер на основі спільних карткових ідентифікаторів та реквізитів карткотримача. Навчено моделі GraphSAGE та GAT — найкраща окрема модель sage\_128\_3L досягла OOF F1-score = 0,8248.

4. Реалізовано стакінг-ансамбль з 19 базових моделей (5 граф-нейронних та 14 табличних) із мета-моделлю XGBoost та rank-нормалізацією OOF-прогнозів. Підсумкове значення OOF F1-score = 0,8877 при precision = 0,91 та recall = 0,86.

5. Виконано валідацію моделі (5-fold StratifiedKFold, bootstrap довірчі інтервали, ablation-study, leave-one-out аналіз) та розроблено трирівневу стратегію інтеграції у платіжний сервіс із сегментованими порогами класифікації.

Проведено детальний аналіз помилок класифікації. «Тихий фрод» (пропущені шахраї) переважно складається з ізольованих вузлів без мережеских

зв'язків, що вказує на напрямок подальших досліджень — розширення графу за рахунок додаткових типів зв'язків (IP, device, email).

Розроблено трирівневу стратегію бізнес-інтеграції (автоблокування з сегментованими порогами (зв'язані:  $\text{score} \geq 0.397$ , ізольовані:  $\text{score} \geq 0.464$ ), ручна перевірка для пограничних випадків, посилений моніторинг для  $\text{score} < 0.20$ ) та систему моніторингу concept drift на основі PSI.

Практична значимість: розроблена система виявляє 86% шахрайських акаунтів з точністю 91% ( $\text{recall}=0.86$ ,  $\text{precision}=0.91$ ), що у масштабах платіжного сервісу з мільйонами користувачів означає потенційну економію десятків мільйонів гривень щорічно. Методологія є універсальною та може бути адаптована для інших фінансових сервісів з мережевими взаємодіями.

Напрямки подальших досліджень визначаються як обмеженнями поточного підходу, так і новими можливостями, що відкриваються з розвитком технологій:

1) Динамічні граф-нейронні мережі (Temporal GNN). Поточний підхід використовує статичний card-sharing граф, побудований на всіх доступних транзакціях. Однак реальна мережа шахрайства еволюціонує: нові зв'язки утворюються, старі розриваються. Temporal GNN, такі як TGN (Temporal Graph Network) та TGAT (Temporal Graph Attention Network), можуть моделювати цю динаміку, виявляючи нові шахрайські кільця на ранній стадії їх формування. Очікуване покращення F1: +0.01-0.03.

2) Інтеграція attention-механізмів. Заміна mean-агрегації у GraphSAGE на attention-зважену агрегацію (GAT — Graph Attention Network) дозволить моделі автоматично визначати, які сусіди є найбільш інформативними для класифікації кожного вузла. Наприклад, сусід з  $\text{fraud\_score}=0.99$  повинен мати більшу вагу, ніж сусід з  $\text{fraud\_score}=0.01$ . Попередні дослідження показують покращення на 0.5-2% для задач класифікації вузлів.

3) Federated Learning для кросбанківського обміну fraud-сигналами. Шахраї часто діють одночасно в кількох банках, але банки не можуть ділитися даними

клієнтів через GDPR та конфіденційність. Federated Learning дозволяє навчати спільну GNN-модель, де кожен банк тренує модель на своїх даних та ділиться тільки градієнтами (не даними). Це може значно покращити виявлення крос-банківського шахрайства без порушення конфіденційності.

4) Використання Large Language Models (LLM) для генерації пояснень рішень моделі. Замість технічних SHAP-значень, LLM може генерувати людиночитабельні пояснення: «Акаунт заблоковано, оскільки 8 з 10 пов'язаних акаунтів раніше були ідентифіковані як шахрайські, а 12 транзакцій відхилено банком-емітентом з причини fraud». Такі пояснення значно покращують ефективність ручної перевірки аналітиками.

5) Self-supervised pre-training для GNN. Навчання GNN на задачах самонавчання (prediction of masked node features, link prediction, graph contrastive learning) перед fine-tuning на задачі класифікації може покращити якість представлень вузлів, особливо для ізольованих або слабозв'язаних вузлів.

6) Мультиграф з гетерогенними типами ребер. Розширення card-sharing графу додатковими типами зв'язків: спільний device fingerprint, IP-адреса, email domain, phone prefix. Гетерогенний GNN (HAN, R-GCN) може навчити різні ваги для різних типів зв'язків, що є більш реалістичним представленням мережі взаємодій.

Загалом, результати дослідження підтверджують, що комбінація граф-нейронних мереж з табличними методами через стакінг є потужним та перспективним підходом до виявлення фінансового шахрайства. Мережеві зв'язки між користувачами є найбільш інформативним сигналом, що відкриває нові можливості для побудови ефективних антифрод-систем нового покоління.

## СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Hamilton W., Ying R., Leskovec J. Inductive Representation Learning on Large Graphs // Advances in Neural Information Processing Systems (NeurIPS). — 2017. — P. 1024–1034.
2. Kipf T.N., Welling M. Semi-Supervised Classification with Graph Convolutional Networks // International Conference on Learning Representations (ICLR). — 2017.
3. Lin T.-Y., Goyal P., Girshick R., He K., Dollár P. Focal Loss for Dense Object Detection // IEEE International Conference on Computer Vision (ICCV). — 2017. — P. 2980–2988.
4. Ke G., Meng Q., Finley T. et al. LightGBM: A Highly Efficient Gradient Boosting Decision Tree // Advances in Neural Information Processing Systems (NeurIPS). — 2017. — P. 3146–3154.
5. Chen T., Guestrin C. XGBoost: A Scalable Tree Boosting System // Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. — 2016. — P. 785–794.
6. Wolpert D.H. Stacked Generalization // Neural Networks. — 1992. — Vol. 5, No. 2. — P. 241–259.
7. Breiman L. Random Forests // Machine Learning. — 2001. — Vol. 45, No. 1. — P. 5–32.
8. Vaswani A. et al. Attention Is All You Need // Advances in Neural Information Processing Systems (NeurIPS). — 2017. — P. 5998–6008.
9. Hochreiter S., Schmidhuber J. Long Short-Term Memory // Neural Computation. — 1997. — Vol. 9, No. 8. — P. 1735–1780.
10. Fey M., Lenssen J.E. Fast Graph Representation Learning with PyTorch Geometric // ICLR Workshop on Representation Learning on Graphs and Manifolds. — 2019.
11. Chawla N.V., Bowyer K.W., Hall L.O., Kegelmeyer W.P. SMOTE: Synthetic Minority Over-sampling Technique // Journal of Artificial Intelligence Research. — 2002. — Vol. 16. — P. 321–357.
12. Bolton R.J., Hand D.J. Statistical Fraud Detection: A Review // Statistical Science. — 2002. — Vol. 17, No. 3. — P. 235–255.
13. Phua C., Lee V., Smith K., Gayler R. A Comprehensive Survey of Data Mining-Based Fraud Detection Research // arXiv preprint. — 2010. — arXiv:1009.6119.
14. Pourhabibi T. et al. Fraud Detection: A Systematic Literature Review of Graph-Based Anomaly Detection Approaches // Decision Support Systems. — 2020. — Vol. 133. — P. 113303.
15. Liu Z., Chen C., Yang X. et al. Heterogeneous Graph Neural Networks for Malicious Account Detection // Proceedings of the 27th ACM International Conference on Information and Knowledge Management. — 2018. — P. 2077–2085.

16. Weber M. et al. Anti-Money Laundering in Bitcoin: Experimenting with Graph Convolutional Networks for Financial Forensics // KDD Workshop on Anomaly Detection in Finance. — 2019.
17. Goodfellow I., Bengio Y., Courville A. Deep Learning. — MIT Press, 2016. — 800 p.
18. Hastie T., Tibshirani R., Friedman J. The Elements of Statistical Learning. — 2nd ed. — Springer, 2009. — 745 p.
19. Nilson Report. Card Fraud Losses Reach \$33.45 Billion // Issue 1234. — 2024.
20. European Parliament and Council. General Data Protection Regulation (GDPR) // Regulation (EU) 2016/679. — 2016.
21. Закон України «Про захист персональних даних» від 01.06.2010 № 2297-VI.
22. Закон України «Про платіжні послуги» від 30.06.2021 № 1591-IX.
23. Grinsztajn L., Oyallon E., Varoquaux G. Why do tree-based models still outperform deep learning on tabular data? // Advances in Neural Information Processing Systems (NeurIPS). — 2022.
24. Xu K., Hu W., Leskovec J., Jegelka S. How Powerful are Graph Neural Networks? // International Conference on Learning Representations (ICLR). — 2019.
25. Ying R., Bourgeois D., You J., Zitnik M., Leskovec J. GNNExplainer: Generating Explanations for Graph Neural Networks // Advances in Neural Information Processing Systems (NeurIPS). — 2019.
26. Lundberg S.M., Lee S.-I. A Unified Approach to Interpreting Model Predictions // Advances in Neural Information Processing Systems (NeurIPS). — 2017. — P. 4765–4774.
27. Arik S.Ö., Pfister T. TabNet: Attentive Interpretable Tabular Learning // Proceedings of the AAAI Conference on Artificial Intelligence. — 2021. — Vol. 35. — P. 6679–6687.
28. Velickovic P., Cucurull G., Casanova A. et al. Graph Attention Networks // International Conference on Learning Representations (ICLR). — 2018.
29. Rashid K.M., Louis J. Window-Warping: A Time Series Data Augmentation of IMU Data for Construction Equipment Activity Identification // Proceedings of ISARC. — 2019.
30. Директива (ЄС) 2015/2366 (PSD2) про платіжні послуги на внутрішньому ринку. Європейський Парламент та Рада. — 2015.
31. Wang J., Wen R., Wu C., Huang Y., Xiong J. FdGars: Fraudster Detection via Graph Convolutional Networks in Online App Review System // ACM SIGKDD. — 2019. — Pp. 310–316.
32. Liu Z., Dou Y., Yu P. S., Deng Y., Peng H. Alleviating the Inconsistency Problem of Applying Graph Neural Network to Fraud Detection // ACM SIGIR. — 2020. — Pp. 1569–1572.

33. Rossi E., Chamberlain B., Frasca F., Eynard D., Monti F., Bronstein M. Temporal Graph Networks for Deep Learning on Dynamic Graphs // ICML Workshop on Graph Representation Learning. — 2020.
34. Liu Y., Jin M., Pan S., Zhou C., Zheng Y., Xia F., Yu P. S. Graph Self-Supervised Learning: A Survey // IEEE Transactions on Knowledge and Data Engineering. — 2022.
35. Chen T., Guestrin C. XGBoost: A Scalable Tree Boosting System // ACM SIGKDD. — 2016. — Pp. 785–794.
36. Veličković P., Cucurull G., Casanova A., Romero A., Liò P., Bengio Y. Graph Attention Networks // International Conference on Learning Representations (ICLR). — 2018.
37. Suprunenko M., Myroshnychenko I. Photorealistic Synthetic Data Generation for Robust Ukrainian Fiscal Receipt OCR : тези доп. наук.-практ. конф. — Київ : КНУ ім. Т. Шевченка, 2026.
38. Suprunenko M., Myroshnychenko I. Information Analytics for Forming and Supporting Productive Habits : тези доп. наук.-практ. конф. — Київ : КНУ ім. Т. Шевченка, 2026.

## ДОДАТКИ

### Додаток А. Конфігурації базових моделей

**Таблиця А.1. Гіперпараметри LightGBM моделей**

Параметр	focal (base)	lgb_deep	lgb_shallow	lgb_dart
num_leaves	255	512	31	255
max_depth	8	12	5	8
learning_rate	0.05	0.03	0.05	0.05
n_estimators	2000	3000	1500	2000
min_child_samples	50	30	100	50
feature_fraction	0.8	0.7	0.9	0.8
bagging_fraction	0.8	0.8	0.9	0.8
boosting_type	gbdt	gbdt	gbdt	dart
Focal alpha	0.75	—	—	—
Focal gamma	2.0	—	—	—
scale_pos_weight	—	25	25	25
OOF F1	0.7109	0.6987	0.6854	0.7021

**Таблиця А.2. Гіперпараметри XGBoost моделей**

Параметр	xgb_v11	xgb_focal	xgb_focal_a80
max_depth	8	8	8
learning_rate	0.05	0.05	0.05
n_estimators	2000	2000	2000
min_child_weight	50	50	50
subsample	0.8	0.8	0.8
colsample_bytree	0.8	0.8	0.8
Focal alpha	—	0.75	0.80
Focal gamma	—	2.0	2.0
scale_pos_weight	25	—	—
OOF F1	0.6993	0.6701	0.6698

## Додаток Б. Архітектура GraphSAGE

**Таблиця Б.1. Конфігурації GraphSAGE моделей**

Параметр	sage_128_3L	sage_64_2L
Hidden dim	128	64
Num layers	3	2
Dropout	0.3	0.2
Learning rate	0.001	0.002
Weight decay	1e-4	1e-4
Epochs (max)	150	150
Batch size	4096	4096
Num neighbors	[15, 10]	[15, 10]
pos_weight	25.5	25.5
Scheduler	CosineAnnealing	CosineAnnealing
Early stop patience	6 (×5 epochs)	6 (×5 epochs)
OOF F1	0.8248	0.8132
OOF AUC	0.9931	0.9925
Training time	280 min	70 min

## Додаток В. Оптимальний ансамбль

Таблиця В.1. Склад фінального стакінг-ансамблю

Модель	Тип	Індивід. OOF F1	Роль в ансамблі
sage_128_3L	GraphSAGE	0.8248	Основа: мережевий сигнал
sage_64_2L	GraphSAGE	0.8132	Різноманітність GNN
focal	LightGBM + FL	0.7109	Найкращий табличний
twostage	LGB 2-stage	0.7106	Двоетапний підхід
rf	Random Forest	0.6040	Bagging різноманітність
lgb_v3	LightGBM	0.7066	Velocity features
xgb_v11	XGBoost	0.6993	Depth-wise бустінг
mlp	Neural Network	0.6353	Нелінійні межі

Мета-модель: LogisticRegression(XGB meta-learner, solver=lbfgs,  
max\_iter=2000)

Нормалізація: rank normalization (rankdata(a) / len(a))

Крос-валідація: StratifiedKFold( $n_{splits}=5$ , shuffle=True, random\_state=42)

Фінальний OOF F1-score: 0.8877

## Додаток Г. Детальні результати крос-валідації

**Таблиця Г.1. OOF F1-score всіх базових моделей**

#	Модель	Тип	OOF F1	OOF AUC
1	sage_128_3L	GraphSAGE	0.8248	0.9931
2	sage_64_2L	GraphSAGE	0.8132	0.9925
3	focal	LGB+FL	0.7109	0.9412
4	twostage	LGB 2- stage	0.7106	0.9398
5	lgb_v9	LightGBM	0.7066	0.9387
6	lgb_dart	LGB DART	0.7021	0.9356
7	xgb_v11	XGBoost	0.6993	0.9345
8	focal_a80	LGB+FL	0.6981	0.9334
9	focal_g2.5	LGB+FL	0.6917	0.9312
10	focal_g1.0	LGB+FL	0.6903	0.9298
11	lgb_deep	LightGBM	0.6987	0.9341
12	lgb_shallow	LightGBM	0.6854	0.9267
13	focal_g3.0	LGB+FL	0.6851	0.9259
14	focal_g4.0	LGB+FL	0.6789	0.9213
15	xgb_focal	XGB+FL	0.6701	0.9187
16	xgb_focal_a80	XGB+FL	0.6698	0.9184
17	xgb_focal_g1.5	XGB+FL	0.6648	0.9156
18	lgb_v3	LightGBM	0.6470	0.9023
19	mlp	MLP	0.6353	0.8945
20	rf	Random Forest	0.6040	0.8756

Продовження таблиці Г.1.

#	Модель	Тип	OOF F1	OOF AUC
21	lgb_v2	LightGBM	0.5960	0.8634
22	extra_trees	Extra Trees	0.5775	0.8523
23	lgb_highrecall	LightGBM	0.5594	0.8412
24	lgb_seed43	LightGBM	0.7098	0.9406
25	lgb_seed44	LightGBM	0.7091	0.9401
26	xgb_seed43	XGBoost	0.6981	0.9338
27	lgb_susp_v1	LGB susp	0.6899	0.9287
28	xgb_susp	XGB susp	0.6823	0.9245
29	lgb_anomaly	LGB anomaly	0.6512	0.9089
30	lstm	LSTM	0.4666	0.7234
31	transformer	Transformer	0.4703	0.7289

*Примітка: моделі відсортовані за F1-score у межах кожного блоку. OOF AUC обчислено на rank-нормалізованих предиктах. Всі моделі навчені з StratifiedKFold(5, shuffle=True, random\_state=42).*

## Додаток Д. Кореляційна матриця OOF-предиктів

Кореляційна матриця Пірсона між OOF-предиктами базових моделей дозволяє оцінити різноманітність ансамблю. Низька кореляція між моделями вказує на їх комплементарність, що є ключовою умовою ефективного стакингу.

**Таблиця Д.1. Кореляція між предиктами ключових моделей**

	sage128	sage64	focal	twostg	rf	lgb_v3	xgb	mlp
sage128	1.00	0.95	0.72	0.71	0.65	0.68	0.70	0.63
sage64	0.95	1.00	0.70	0.69	0.63	0.66	0.68	0.61
focal	0.72	0.70	1.00	0.94	0.82	0.89	0.92	0.78
twostg	0.71	0.69	0.94	1.00	0.81	0.88	0.91	0.77
rf	0.65	0.63	0.82	0.81	1.00	0.78	0.80	0.73
lgb_v3	0.68	0.66	0.89	0.88	0.78	1.00	0.87	0.75
xgb	0.70	0.68	0.92	0.91	0.80	0.87	1.00	0.76
mlp	0.63	0.61	0.78	0.77	0.73	0.75	0.76	1.00

Аналіз кореляційної матриці виявляє чітку блочну структуру:

1) GNN-блок (sage128, sage64): висока внутрішня кореляція (0.95), низька кореляція з табличними моделями (0.61-0.72). GNN використовує графову структуру, табличні — тільки ознаки вузла.

2) Табличний GBM-блок (focal, twostage, lgb\_v3, xgb): висока внутрішня кореляція (0.87-0.94), оскільки всі використовують gradient boosting на схожих ознаках.

3) RF та MLP мають нижчу кореляцію з GBM-блоком (0.73-0.82), що підтверджує їх корисність для різноманітності ансамблю.

4) Найнижча кореляція — між GNN та MLP (0.61-0.63), що робить їх найбільш комплементарною парою.