

Міністерство освіти і науки України
Національний технічний університет
«Дніпровська політехніка»

Факультет інформаційних технологій
(факультет)

Кафедра системного аналізу та управління
(повна назва)

ПОЯСНЮВАЛЬНА ЗАПИСКА
кваліфікаційної роботи ступеня магістра

Студента Жучкова Сергія Ігоровича
академічної групи 124М-23-1
спеціальності 124 Системний аналіз

на тему: Дослідження та порівняння ефективності методів класифікації даних у
задачах Fraud Detection»

Керівники	Прізвище, ініціали	Оцінка за шкалою		Підпис
		рейтингово ю	Інституційною	
кваліфікаційної роботи				
розділів:				
Інформаційно- аналітичний				
Спеціальний розділ				
Рецензент				
Нормоконтролер				

Дніпро
2024

ЗАТВЕРДЖЕНО:
завідувач кафедри
Системного аналізу та управління
(повна назва)

_____ к.т.н., доц. Желдак Т.А.
(підпис) (прізвище, ініціали)

« _____ » _____ 20__ року

ЗАВДАННЯ
на кваліфікаційну роботу
ступеня магістра

студенту Жучкову С.І. академічної групи 124М- 23-1
спеціальності: 124 Системний аналіз
на тему «Дослідження та порівняння ефективності методів класифікації даних
у задачах Fraud Detection»
затверджену наказом ректора НТУ «Дніпровська політехніка»
від 16.10.2024 р. №1388 – С

Розділ	Зміст	Терміни виконання
1. Інформаційно-аналітичний розділ	<i>Визначити предметну область дослідження та проблему, що розв'язується. Обґрунтувати методи виконання поставлених завдань.</i>	10.09.2024 – 01.11.2024
2. Спеціальний розділ	<i>Розв'язати поставлені задачі: описати структуру об'єкта досліджень, застосувати методи класифікації та провести експерименти для підвищення якості моделей.</i>	01.11.2024 – 30.11.2024

Завдання видано _____ (підпис) _____ (прізвище, ініціали)

Дата видачі: _____

Дата подання до екзаменаційної комісії: _____

Прийнято до виконання _____ (підпис студента)

Жучков С.І.
(прізвище, ініціали)

РЕФЕРАТ

Пояснювальна записка: 138 с., 59 рис., 6 табл., 6 додатків, 20 джерел.

Ключові слова: МОДЕЛЬ, КЛАСИФІКАЦІЯ, НЕЗБАЛАНСОВАНИЙ ДАТАСЕТ, МЕТОДИ БАЛАНСУВАННЯ ДАНИХ, КЛАСТЕРИЗАЦІЯ, ШАХРАЙСТВО, ЯКІСТЬ, КЛАСТЕР, КЛАСТЕРИЗАЦІЯ, ПРОГНОЗУВАННЯ, EXBOOST, LIGHTGBM, CATBOOST.

Об'єктом дослідження кваліфікаційної роботи є аналіз нових методів прогнозування шахрайських транзакцій, де основним аспектом є вплив незбалансованості класів у наборах даних на якість моделей прогнозування.

Предметом дослідження є методи класифікації незбалансованих даних, а саме їх теоретичні аспекти та практичне застосування на реально існуючих даних і оцінка їх ефективності.

Метою кваліфікаційної роботи є дослідження, застосування та порівняння методів класифікації незбалансованих наборів даних в задачах виявлення фроду, також знаходження підходящих методів кластеризації для покращення якості моделі прогнозування.

Результати кваліфікаційної роботи включають детальний аналіз та оцінку методів класифікації незбалансованих даних, включаючи їх ефективність прогнозування шахрайських транзакцій.

Використовується мова програмування Python із бібліотеками, такими як Pandas, NumPy, Matplotlib, Seaborn та Scikit-learn, для обробки даних, виконання аналізу та прогнозування. Інноваційність роботи відрізняється своїм підходом до вибору та застосування методів, які оптимально підходять для обробки незбалансованих даних.

Ця робота вписується у ширший контекст досліджень у сфері машинного навчання та аналізу даних, зокрема у дослідження, що стосуються класифікації незбалансованих даних.

ABSTRACT

Explanatory note: 138 pages, 59 figures, 6 tables, 6 appendices, 20 references.

Keywords: MODEL, CLASSIFICATION, IMBALANCED DATASET, DATA BALANCING METHODS, CLUSTERING, FRAUD, QUALITY, CLUSTER, CLUSTERING, PREDICTION, XGBOOST, LIGHTGBM, CATBOOST.

The object of this qualification work is the analysis of new methods for predicting fraudulent transactions, with a primary focus on the impact of class imbalance in datasets on the performance of predictive models.

The subject of the study is the classification methods for imbalanced data, specifically their theoretical aspects and practical application to real-world data, along with the evaluation of their effectiveness.

The aim of the qualification work is to explore, apply, and compare classification methods for imbalanced datasets in fraud detection tasks, as well as to identify suitable clustering methods to improve the quality of predictive models.

The results of the qualification work include a detailed analysis and evaluation of classification methods for imbalanced data, highlighting their effectiveness in predicting fraudulent transactions.

The study employs Python programming language with libraries such as Pandas, NumPy, Matplotlib, Seaborn, and Scikit-learn for data processing, analysis, and prediction. The innovative aspect of this work lies in its approach to selecting and applying methods optimally suited for handling imbalanced data.

This work contributes to the broader field of machine learning and data analysis, specifically in research related to the classification of imbalanced datasets.

ЗМІСТ

ВСТУП	6
1 ІНФОРМАЦІЙНО-АНАЛІТИЧНИЙ РОЗДІЛ	9
1.1 Дослідження банківських транзакцій та особливості датасету	9
1.2 Методи підвищення якості прогнозу моделей класифікації	10
1.3 Нерівномірний розподіл класів та проблема нерівномірної класифікації	16
1.4 Методи балансування даних	17
1.5 Методи групування даних	25
1.6 Методи класифікації незбалансованих наборів даних	29
1.7 Методи оптимізації параметрів моделей	38
1.8 Показники ефективності моделей класифікації	40
Висновок	47
2 СПЕЦІАЛЬНИЙ РОЗДІЛ	48
2.1 Постановка задачі	48
2.1.1 Мета дослідження	48
2.1.2 Задачі дослідження	48
2.2 Обґрунтування вибору середовища програмування	49
2.2 Exploratory Data Analysis	52
2.4 Побудова моделей прогнозування	75
2.4.1 Модель прогнозування з кластеризацією K-means	75
2.4.2 Метод Fuzzy Clustering	90
2.4.3 Метод CatBoost	97
2.4.4 Метод LightGBM	103
ВИСНОВКИ	109
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ	111
ДОДАТОК А. Відомість матеріалів кваліфікаційної роботи	114
ДОДАТОК Б	115
ДОДАТОК В	116
ДОДАТОК Г	125
ДОДАТОК Д	131
ДОДАТОК Е	134

ВСТУП

Останніми роками фірми, які надають фінансові послуги, звертаються до технологій машинного навчання, щоб допомогти виявити та, за результатами машинного навчання, запобігти шахрайству. Ця технологія машиного навчання, яка використовує алгоритми для виявлення шаблонів у великих наборах даних, є ефективним інструментом для виявлення та запобігання шахрайству.

Галузь фінансових послуг особливо вразлива для шахраїв, оскільки вони мають доступ до величезних обсягів даних клієнтів. У результаті традиційних методів виявлення шахрайства часто недостатньо для захисту клієнтів та їхніх коштів. Машинне навчання швидко стає важливим інструментом для запобігання шахрайству та захисту активів клієнтів.

Алгоритми машинного навчання можуть виявляти шаблони в даних, які люди можуть не бачити. Це полегшує ідентифікацію підозрілої активності, наприклад підозріло великих сум транзакцій або кількох транзакцій з одного облікового запису.

Аналізуючи дані клієнтів у режимі реального часу, банки можуть швидко виявляти та запобігати шахрайству та таким чином допомагає захистити клієнтів від втрат і зберігає безпеку фінансової системи.

Мета кваліфікаційної роботи полягає у застосуванні та порівнянні методів класифікації незбалансованих наборів даних в задачах виявлення фроду.

Предметом дослідження є методи класифікації незбалансованих даних, а саме їх теоретичні аспекти та практичне застосування до реальних даних включно із оцінкою їх ефективності.

Об'єктом дослідження є процес класифікації незбалансованих наборів даних в задачах виявлення фроду.

Для досягнення мети дослідження потрібно вирішити такі задачі:

1. Провести аналіз теоретичних аспектів класифікації даних, зокрема методів роботи з незбалансованими даними.

2. Здійснити огляд методів класифікації незбалансованих наборів даних, включаючи методи, засновані на деревах рішень, нейронних мережах, ансамблевих методах та інших.

3. Провести експерименти з використанням методів класифікації на даних з задачі виявлення фроду.

4. Здійснити аналіз результатів експерименту та сформулювати висновки щодо ефективності кожного методу.

Теоретичне значення кваліфікаційної роботи магістра спрямована на розробку теоретичних засад і практичних підходів до класифікації незбалансованих даних у задачах виявлення шахрайства. У ході дослідження буде сформульовано наукові положення, які розширюють розуміння механізмів адаптації класифікаційних алгоритмів до обробки асиметричних наборів даних, зокрема шляхом оптимізації метрик, використання спеціалізованих алгоритмів та підходів до ресемплінгу даних. Отримані результати можуть бути використані для вдосконалення існуючих методів аналізу даних, що матиме значення для галузей, пов'язаних із великими обсягами транзакційних даних.

Наукова новизна отриманих результатів. Вперше одержано комплексну оцінку ефективності популярних методів класифікації незбалансованих даних (XGBoost, CatBoost, LightGBM тощо) з урахуванням специфіки задачі виявлення шахрайства. Вдосконалено методи підвищення точності класифікації за рахунок оптимізації порогових значень для мінімізації помилок класифікації у класах із низькою частотою. Дістало подальший розвиток підходи до використання алгоритмів ресемплінгу, таких як SMOTE, з інтеграцією у процес навчання моделей для покращення виявлення рідкісних класів. Кожен із зазначених результатів чітко формулює сутність запропонованих рішень та демонструє їхню наукову новизну, забезпечуючи легке сприйняття та можливість практичного застосування.

Практичне значення отриманих результатів. Розроблені методи дозволяють підвищити точність і надійність автоматизованих систем виявлення шахрайства у фінансових установах, що сприятиме зменшенню фінансових втрат

клієнтів та банків. Результати роботи можуть бути інтегровані у програмне забезпечення для моніторингу транзакцій, забезпечуючи адаптацію до нових шаблонів шахрайства. Рекомендації щодо застосування методів класифікації та роботи з незбалансованими даними є універсальними і можуть бути використані у суміжних сферах, наприклад, у задачах виявлення аномалій у кібербезпеці або моніторингу поведінки клієнтів.

1 ІНФОРМАЦІЙНО-АНАЛІТИЧНИЙ РОЗДІЛ

1.1 Дослідження банківських транзакцій та особливості датасету

Дослідження банківських транзакцій є важливою складовою фінансового аналізу. У зв'язку зі зростанням обсягів банківських транзакцій, з'явилася потреба в розробці ефективних методів аналізу та виявлення аномальних транзакцій, які можуть вказувати на шахрайську діяльність. Такі датасети зазвичай містять інформацію про час, коли була проведена транзакція, типу транзакції, сумі операції в національній валюті, ім'я клієнта, що почав операцію, початковий баланс до операції, баланс клієнта та одержувача після транзакції, ідентифікатор одержувача транзакції та чи є транзакція шахрайством або ні.

Одразу можна зрозуміти головну особливість таких датасетів – категоріальні дані (тип та ідентифікатори клієнта та одержувача). Так як моделі класифікації можуть вчитися лише на числових даних, ці ознаки, якщо з ними не розібратися, ускладнять роботу з даними для цих моделей, тому ці ознаки потребують методів кодування категоріальних ознак, такі як Label Encoding або One-Hot Encoding.

Для успішного та ефективного розуміння даних банківських транзакцій потрібно провести EDA (розвідковий аналіз даних, англ. exploratory data analysis) та залучити релевантні методи кластеризації.

Exploratory Data Analysis – процес дослідження даних з метою зрозуміти їх структуру, виявити закономірності та корисну інформацію. Він складається з наступних етапів:

1. Первинний аналіз даних: ознайомлення з даними, їх структурою та типами, виявлення пропущених значень, викинутих стовпців і рядків, статистичний аналіз розподілу даних.

2. Візуалізація даних: відображення даних у вигляді графіків, діаграм, хмарних графіків, гістограм, що дозволяє отримати зображення розподілу даних та співвідношення між ними.

3. Розуміння залежностей: пошук залежностей між різними змінними за допомогою кореляційних матриць, теплових карт, факторного аналізу та інших методів.

4. Обробка даних: виявлення аномальних даних, опрацювання пропущених значень, вибір змінних для подальшого аналізу та побудови моделей.

5. Перевірка гіпотез: тестування гіпотез щодо залежності між змінними, відповідності даних певному розподілу, перевірка статистичних гіпотез.

6. Висновки: на основі проведеного аналізу формулювання висновків та рекомендацій щодо подальшого використання даних, побудови моделей або виконання подальших досліджень.

EDA є важливим етапом у будь-якому проекті з аналізу даних, оскільки він дозволяє зрозуміти характеристики та особливості даних, з якими ми маємо справу. Це може допомогти визначити, які моделі, алгоритми та методи аналізу можуть бути найбільш ефективними для нашої задачі, а також дозволяє виявити потенційні проблеми та виклики, пов'язані з даними.

1.2 Методи підвищення якості прогнозу моделей класифікації

Методи підвищення точності моделей класифікації включають різні техніки, які можуть бути використані для покращення якості прогнозів моделі.

Розглянемо деякі з них:

1. Оптимізація даних:

Очищення та підготовка даних є важливим етапом у будь-якому проекті машинного навчання. Неочищені або погано підготовлені дані можуть значно вплинути на якість моделі та її точність.

На етапі очищення даних виконують перевірку на наявність відсутніх значень та їх обробку (видалення рядків з відсутніми значеннями, заповнення пропущених даних тощо) та видалення дублікатів.

Видалення викидів:

Викиди (англ. outliers)-це значення, що виділяється із звичайного діапазону даних і може бути помилкою вимірювання або відображенням особливостей вихідних даних. Наявність викиди у даних може негативно вплинути на точність моделі машинного навчання, оскільки модель може зосередитись на неправильних залежностях і, отже, недооцінити або переоцінити результат.

Існують різні методи обробки викидів, зокрема, видалення або заміна викидів. У випадку видалення викидів, можна використовувати статистичні методи для визначення діапазону "нормальних" значень, а потім видалити ті значення, що виходять за межі цього діапазону. Інший підхід - заміна значень, що виходять за межі діапазону, на інші значення, наприклад, на середнє або медіану даних.

Однак, не завжди видалення або заміна викидів є корисним підходом. Наприклад, у деяких випадках, викид може бути корисним сигналом для моделі, особливо, якщо він представляє реальну аномалію, яка може мати значення для вирішення задачі. Тому перед застосуванням методів обробки викидів, необхідно ретельно проаналізувати дані та зважити на вплив викидів на точність моделі.

Encoding або кодування категоріальних даних:

Кодування категоріальних даних являє собою процес перетворення категоріальних змінних, які приймають обмежену кількість можливих значень, на числові значення, які можуть бути використані в моделях машинного навчання. Це необхідно, оскільки більшість алгоритмів машинного навчання не можуть працювати з категоріальними змінними у вигляді текстових міток.

Існує кілька способів кодування категоріальних даних, включаючи:

- Бінарне кодування (One-Hot Encoding): кожне можливе значення категорії представлено як бінарний вектор. Наприклад, якщо ми маємо змінну "кольори" з трьома можливими значеннями-червоний, зелений та синій, то ми можемо створити три нові змінні "колір_червоний", "колір_зелений" та "колір_синій", і призначити їм значення 1 або 0, в залежності від того, який колір є для кожного запису. Цей метод використовується для збереження всієї інформації про категорію, але призводить до збільшення кількості ознак у даних.

- Частотне кодування (Frequency Encoding): кожне значення категорії замінюється його частотою в даних. Наприклад, якщо ми маємо змінну "кольори" з трьома можливими значеннями - червоний, зелений та синій, то ми можемо замінити кожен колір на його частоту, тобто кількість разів, коли він зустрічається в даних. Цей метод простіший, ніж бінарне кодування, але може викликати проблеми, якщо значення категорії мають однакову частоту, оскільки вони будуть мати однакове закодоване значення.

- Label Encoding: кожна категорія замінюється на числове значення, яке відповідає її позиції в послідовності унікальних значень. Наприклад, якщо ми маємо змінну "кольори" з трьома можливими значеннями - червоний, зелений та синій, то ми можемо замінити кожен колір на числа 1, 2 та 3 відповідно. Цей метод використовується, коли порядок категорій не має значення, але може призвести до помилкового розуміння, що значення категорій мають певний порядок.

- Target Encoding: кожне значення категорії замінюється на середнє значення цільової змінної для цієї категорії. Наприклад, якщо ми маємо змінну "місто" та цільову змінну "ціна", то ми можемо замінити кожне місто на середню ціну в цьому місті. Цей метод використовується, коли категорії можуть мати різний вплив на цільову змінну, але може призвести до перенавчання моделі, якщо певні категорії мають надмірний вплив на цільову змінну.

Масштабування або нормалізація числових змінних:

Масштабування або нормалізація числових змінних-процес приведення значень числових змінних до певного діапазону або розподілу. Це робиться для того, щоб зробити числові змінні порівнянними, а також для підвищення швидкості та якості навчання моделі.

Масштабування може здійснюватися за допомогою двох підходів: нормалізації та стандартизації. Нормалізація зазвичай використовується, коли діапазон значень змінних різний, а стандартизація - коли змінні мають різні дисперсії.

При нормалізації значення змінної масштабуються до діапазону від 0 до 1 або від -1 до 1. Це робиться шляхом віднімання мінімального значення змінної та поділу на різницю між максимальним та мінімальним значеннями змінної.

При стандартизації значення змінної масштабуються так, щоб вони мали середнє значення 0 та стандартне відхилення 1. Це робиться шляхом віднімання середнього значення змінної та поділу на стандартне відхилення.

Масштабування змінних дозволяє моделі краще узгоджуватися з даними та підвищує її точність та стійкість.

2. Feature selection або вибор ознак:

Відбір ознак-процес вибору підмножини ознак з вхідних даних, які є найбільш інформативними для побудови моделі машинного навчання. Це може бути корисним для зменшення кількості ознак, що використовуються в моделі, що дозволяє зменшити ризик перенавчання, збільшити швидкість навчання та покращити загальну ефективність моделі.

Існує багато методів відбору ознак, такі як:

- Filter methods (методи фільтрації): використовують статистичні метрики для відбору ознак, які мають найбільший вплив на вихідну змінну. Приклади таких метрик - кореляційний коефіцієнт Пірсона, mutual information, chi-squared test тощо.

- Wrapper methods (методи упаковки): використовують ітеративний підхід, де модель машинного навчання навчається на різних підмножинах ознак і відбирається найкраща підмножина на основі валідаційних метрик. Приклади таких методів - Recursive Feature Elimination, Forward Selection, Backward Elimination тощо.

- Embedded methods (вбудовані методи): відбір ознак вбудований у процес навчання моделі машинного навчання. Найбільш відомі приклади вбудованих методів - Lasso Regression, Elastic Net, Ridge Regression тощо.

3. Cross-validation або кросс-валідація:

Кросс-валідація-метод оцінки якості моделі машинного навчання, який дозволяє оцінити, наскільки добре модель буде працювати на нових невідомих

даних. Зазвичай кросс-валідація використовується для оцінки якості моделі перед її застосуванням на реальних даних.

Ідея полягає в тому, що наявний набір даних розділяється на кілька частин, наприклад, 5 частин (5-fold cross validation) або 10 частин (10-fold cross validation). Далі модель навчається на 4 частинах, а потім перевіряється на залишковій 1 частині. Ця процедура повторюється кілька разів, при цьому кожна з частин використовується як тестовий набір даних, а решта - як тренувальний.

Після цього обчислюється середнє значення метрики якості (наприклад, точності або F1-мери) для всіх тестових наборів даних. Ця метрика дозволяє оцінити якість моделі, тобто наскільки добре вона здатна передбачати класи для нових невідомих даних.

Використання кросс-валідації дозволяє уникнути перенавчання моделі на конкретних даних і забезпечити більш точну оцінку її роботи на нових даних.

4. Оптимізація гіперпараметрів моделі:

Оптимізація гіперпараметрів моделі є процесом знаходження оптимальних значень гіперпараметрів моделі, які забезпечують кращу продуктивність моделі на тестових даних.

Зазвичай, для оптимізації гіперпараметрів моделі використовують методи пошуку гіперпараметрів, такі як решітчатий пошук, випадковий пошук або оптимізація градієнтного спуску. У решітчатому пошуку використовуються певні значення гіперпараметрів, а потім обчислюється метрика якості моделі для кожної комбінації значень гіперпараметрів. Випадковий пошук використовує випадкові значення гіперпараметрів для обчислення метрики якості моделі. Оптимізація градієнтного спуску використовує алгоритм градієнтного спуску для знаходження оптимальних значень гіперпараметрів.

Зазвичай, оптимізація гіперпараметрів моделі відбувається шляхом поділу даних на навчальний та тестовий набори, а потім знаходження оптимальних значень гіперпараметрів моделі на навчальному наборі даних з використанням методів пошуку гіперпараметрів. Після знаходження оптимальних значень

гіперпараметрів моделі на навчальному наборі даних, модель тестується на тестовому наборі даних, щоб перевірити її продуктивність.

5. Ансамблеві моделі машинного навчання:

Ансамбль моделей-це підхід до побудови моделі, який поєднує кілька індивідуальних моделей з метою отримання більш точних та надійних прогнозів. Замість використання окремої моделі, ансамбль об'єднує прогнози кількох моделей та робить узагальнений прогноз на основі їх результатів.

Існує кілька типів ансамблів моделей, найпоширеніші з яких такі:

- Беггінг (Bagging): використовується для зменшення варіації моделі шляхом тренування кількох моделей на різних підмножинах даних. Кожна модель незалежно тренується, а результати їх прогнозів комбінуються, наприклад, шляхом голосування або усереднення.

- Випадковий ліс (Random Forest): являє собою тип ансамблю зачеплених моделей, де використовується рішучий ліс як базова модель. Рішучий ліс використовує дерева рішень для прогнозування і комбінує результати декількох дерев для отримання кінцевого прогнозу.

- Послідовний стекінг (Sequential Stacking): використовується для комбінування результатів кількох моделей, де вихідні дані однієї моделі використовуються як вхідні дані для іншої моделі. Результати прогнозів кожної моделі комбінуються для отримання кінцевого прогнозу.

- Бустінг (Boosting): метод, в якому моделі будуються послідовно, при цьому кожна наступна модель спрямовується на коригування помилок попередніх моделей. Бустінг створює сильну модель шляхом комбінування слабких моделей.

Ансамбль моделей може бути підібраний та налаштований для досягнення кращих результатів. Однак, використання ансамблювання може призвести до більш складної моделі, що вимагає більше обчислювальних ресурсів та часу на навчання.

1.3 Нерівномірний розподіл класів та проблема нерівномірної класифікації

Нерівномірність розподілу класів-ситуація, коли дані, які необхідно класифікувати, мають нерівномірний розподіл по класам. Якщо один клас має набагато більше прикладів, ніж інший, то модель буде схильною до прийняття рішень на користь більшого класу, що може призвести до невдалих прогнозів для меншої класу. Проблема нерівномірної класифікації полягає у тому, що модель не може адекватно виконувати класифікацію меншості, і як результат, може допускати помилки в її класифікації.

Для вирішення проблеми нерівномірної класифікації можна використовувати різні методи, такі як:

- Збільшення вибірки меншої класу (oversampling): цей метод полягає в копіюванні прикладів меншої класу з метою збільшення її розміру до рівня більшої класу. Це дозволяє зменшити проблему нерівномірності розподілу класів, але може призвести до перенавчання моделі.

- Зменшення вибірки більшої класу (undersampling): цей метод полягає у випадковому видаленні прикладів більшої класу до досягнення рівномірного розподілу класів. Даний метод також може призвести до перенавчання моделі.

- Використання ваг для класів (class weighting): цей метод полягає в призначенні ваг класам в залежності від їх розміру. Наприклад, меншому класу можна призначити більшу вагу, щоб зменшити ймовірність помилки класифікації цього класу. Цей метод дозволяє зменшити проблему нерівномірності розподілу класів, не призводячи до перенавчання моделі.

- Використання алгоритмів з урахуванням нерівномірної класифікації (class-imbalanced algorithms): деякі алгоритми машинного навчання, такі як XGBoost, LightGBM та CatBoost, мають вбудовані методи для вирішення проблеми нерівномірної класифікації. Ці алгоритми дозволяють встановлювати різні параметри, такі як `scale_pos_weight`, який дозволяє налаштувати ваги класів.

- При вирішенні проблеми нерівномірної класифікації важливо вибрати метод, який дозволяє зменшити нерівномірність розподілу класів, не призводячи до перенавчання моделі і забезпечує якісну класифікацію для всіх класів.

1.4 Методи балансування даних

Методи балансування даних є важливими для вирішення проблеми нерівномірного розподілу класів у наборі даних. Коли класи у наборі даних представлені нерівномірно, що може призвести до поганої узгодженості моделі та неправильних прогнозів.

Основна причина, чому потрібно застосовувати методи балансування даних, полягає в тому, що моделі машинного навчання мають тенденцію віддавати перевагу класу з більшою кількістю прикладів, вважаючи його більш представником та важливим. Це може призвести до неправильної класифікації меншої кількості класу або ігнорування його зовсім.

Методи балансування даних допомагають уникнути цих проблем та покращують якість моделі класифікації. Вони забезпечують більш справедливу репрезентацію обох класів, зменшуючи перекося і покращуючи здатність моделі розрізняти та класифікувати обидва класи.

У цьому підрозділі ми ознайомимось з ними детальніше:

1. Метод *Oversampling* є одним з підходів до балансування даних, що полягає у створенні додаткових прикладів з меншої кількості класу, щоб збалансувати нерівномірність розподілу класів у наборі даних.

Один з підходів до збільшення кількості записів - дублювання вже наявних записів. Однак цей метод може призвести до перенавчання моделі, оскільки він додає ті самі записи знову та знову, збільшуючи ймовірність переобучення.

Інший підхід до *oversampling* полягає в застосуванні методу *SMOTE*, що створює нові записи, шляхом генерації синтетичних прикладів. *SMOTE* вибирає випадкові записи з меншої кількості класу і використовує їх, щоб створити синтетичні записи, додавши до них випадкові значення між оригінальними

записами. SMOTE дозволяє створити нові приклади, що допомагають уникнути проблем перенавчання та збільшують різноманітність даних.

Алгоритм SMOTE полягає в наступних кроках:

1. Вибір випадкового запису з меншістю класу.
2. Визначення його k найближчих сусідів в цьому ж класі за допомогою метрики відстані, наприклад, Евклідової відстані.
3. Вибір випадкового запису з цих k сусідів.
4. Генерація нового запису попередньо не відомого класу на відрізьку, який з'єднує початковий запис та випадкового сусіда з меншості класу, використовуючи формулу:

$$\text{new_record} = \text{original_record} + (\text{random_number} * (\text{neighbor_record} - \text{original_record}))$$

де random_number - випадкове число від 0 до 1, neighbor_record - запис в меншій кількості класу, вибраний на кроці 3.

5. Повторення кроків 1-4 задану кількість разів, щоб отримати більшу кількість штучних записів.

Для кращого розуміння цього методу візуалізуємо його роботу над даними випадкового датасету:

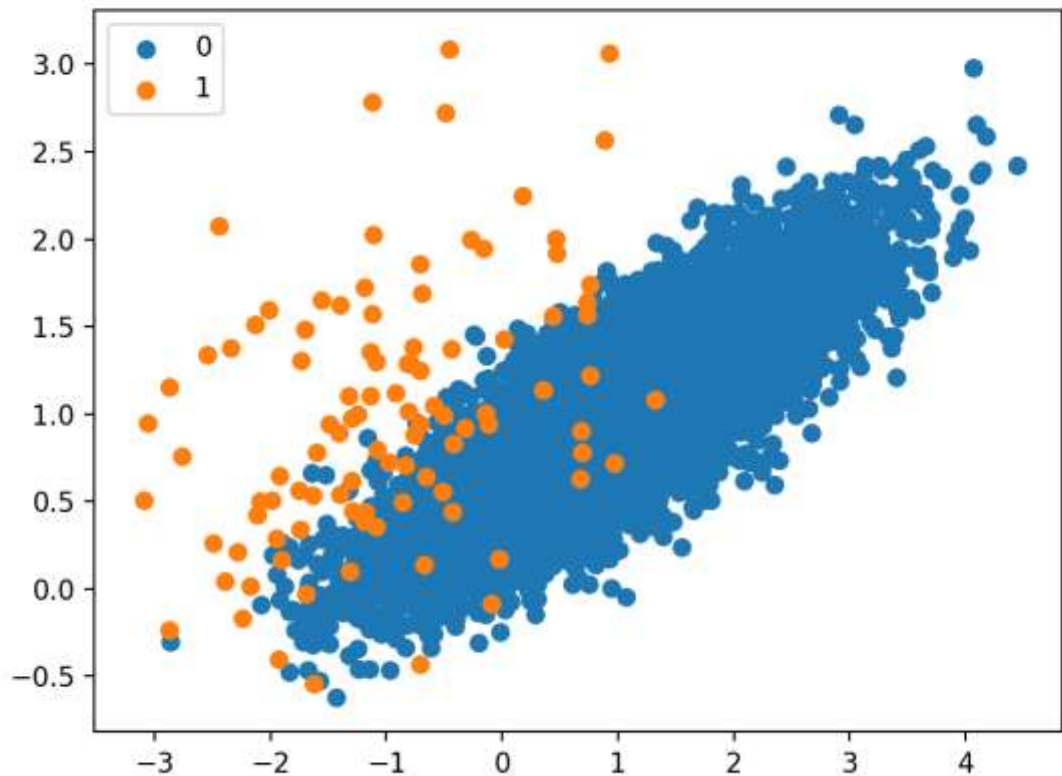


Рисунок 1.1-Графік розділення класу більшості(синім) та клас меншості(помаранчевий)

Ця точкова діаграма набору даних показує велику масу точок, які належать до класу більшості (синій), і невелику кількість точок, розподілених для класу меншості (помаранчевий). Ми бачимо певну міру перекриття між двома класами.

Ми маємо перед собою незбалансований набір, тому за допомогою методу SMOTE створимо синтетичні записи, додавши до них випадкові значення між оригінальними записами.

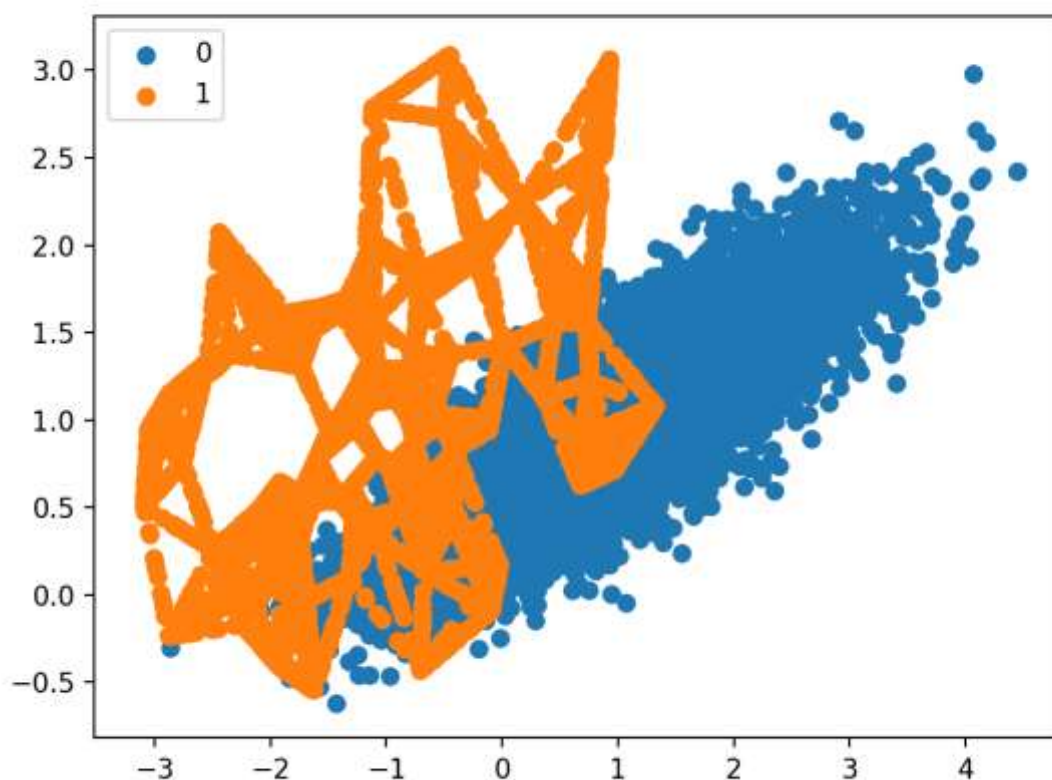


Рисунок 1.2-Результат методу SMOTE (джерело <https://habr.com/ru/articles/842480/>)

Як можемо побачити, метод згенерував нові записи та доповнив їх до класу меншості та тепер графік показує багато інших прикладів у класі меншості, створених уздовж ліній між вихідними прикладами в класі меншості.

Основною перевагою SMOTE є те, що він допомагає уникнути перенавчання моделі, що може статися при повторному використанні наявних записів, що належать до меншості класу. Використання синтетичних записів дозволяє моделі навчитися різним характеристикам меншості класу та забезпечує більш рівномірне розподіл класів.

Однак, SMOTE також має свої недоліки. Наприклад, якщо дані складаються з багатьох дуже вузьких кластерів, то SMOTE може призвести до створення більш широких кластерів та знизити точність моделі. Крім того, SMOTE може не підходити для деяких даних, якщо вони мають велику кількість шуму або випадкових даних, що не пов'язані з конкретним класом. Також, важливо

пам'ятати, що SMOTE не вирішує проблему зі збалансуванням класів, якщо головною причиною їх нерівномірності є неправильна підготовка даних або поганий вибір ознак.

Таким чином, SMOTE є корисним інструментом для розв'язання проблеми нерівномірного розподілу класів у наборі даних, але він має свої обмеження і вимагає обережного застосування.

Algorithm *SMOTE*(T , N , k)
Input: Number of minority class samples T ; Amount of SMOTE $N\%$;
 Number of nearest neighbors k
Output: $(N/100) * T$ synthetic minority class samples

1. (* If N is less than 100%, randomize the minority class samples as only a random percent of them will be SMOTEd. *)
2. **if** $N < 100$
3. **then** Randomize the T minority class samples
4. $T = (N/100) * T$
5. $N = 100$
6. **endif**
7. $N = (int)(N/100)$ (* The amount of SMOTE is assumed to be in integral multiples of 100. *)
8. $k =$ Number of nearest neighbors
9. $numattrs =$ Number of attributes
10. $Sample[] []$: array for original minority class samples
11. $newindex$: keeps a count of number of synthetic samples generated, initialized to 0
12. $Synthetic[] []$: array for synthetic samples
 (* Compute k nearest neighbors for each minority class sample only. *)
13. **for** $i \leftarrow 1$ to T
14. Compute k nearest neighbors for i , and save the indices in the $nnarray$
15. Populate(N , i , $nnarray$)
16. **endfor**

Populate(N , i , $nnarray$) (* Function to generate the synthetic samples. *)

17. **while** $N \neq 0$
18. Choose a random number between 1 and k , call it nn . This step chooses one of the k nearest neighbors of i .
19. **for** $attr \leftarrow 1$ to $numattrs$
20. Compute: $dif = Sample[nnarray[nn]][attr] - Sample[i][attr]$
21. Compute: $gap =$ random number between 0 and 1
22. $Synthetic[newindex][attr] = Sample[i][attr] + gap * dif$
23. **endfor**
24. $newindex++$
25. $N = N - 1$
26. **endwhile**
27. **return** (* End of Populate. *)

End of Pseudo-Code.

Рисунок 1.3-Псевдокод методу SMOTE(джерело [<https://nuancesprog.ru/p/15938/>])

2. Інший метод, що базується на SMOTE, ADASYN (Adaptive Synthetic Sampling), алгоритм, що збільшує кількості зразків меншості у наборі даних, що дозволяє зменшити нерівномірність розподілу класів та враховує щільність класів у просторі ознак.

ADASYN (Adaptive Synthetic Sampling) являє собою алгоритм для збільшення кількості зразків меншості у наборі даних, що дозволяє зменшити нерівномірність розподілу класів. Він працює на основі SMOTE (Synthetic Minority Over-sampling Technique), але враховує щільність класів у просторі ознак.

Основна ідея ADASYN полягає в тому, що він генерує штучні зразки меншості шляхом вибору точки з меншості та її найближчих сусідів. Однак замість генерації нового зразка на прямій, що з'єднує вихідну точку і одного з її сусідів, ADASYN генерує зразки в околі цієї прямої залежно від рівня нерівномірності класів.

Конкретніше, для кожної точки меншості ADASYN обчислює вагу, яка відображає, наскільки багато зразків меншості потрібно синтезувати для цієї точки. Вага розраховується як відношення кількості зразків більшості, які знаходяться в околі точки, до загальної кількості зразків меншості в цьому околі. Очевидно, що чим менше прикладів меншості в околі, тим більша вага надається точці меншості.

Після обчислення ваг для кожної точки меншості, алгоритм вибирає кількість зразків, яку потрібно створити, і для кожної точки вибирає найближчі сусіди з більшості. Зразки генеруються в околі лінії, що з'єднує точку меншості і вибраного сусіда з більшості, пропорційно вагам точок меншості.

Однією з переваг ADASYN є те, що він адаптується до нерівномірного розподілу класів у просторі ознак і створює більше штучних зразків в областях, де даних менше, а менше зразків там, де даних більше. Це дозволяє зменшити ризик перенавчання моделі та покращити її загальну здатність до узагальнення.

Алгоритм ADASYN дуже зрозуміло описав Haibo He, Yang Bai та Edwardo A. Garcia у своїй статті [1] про ADASYN :

Input

(1) Training data set D_{tr} with m samples $\{x_i, y_i\}$, $i = 1, \dots, m$, where x_i is an instance in the n dimensional feature space X and $y_i \in Y = \{1, -1\}$ is the class identity label associated with x_i . Define m_s and m_l as the number of minority class examples and the number of majority class examples, respectively. Therefore, $m_s \leq m_l$ and $m_s + m_l = m$.

Procedure

(1) Calculate the degree of class imbalance:

$$d = m_s/m_l \quad (1)$$

where $d \in (0, 1]$.

(2) If $d < d_{th}$ then (d_{th} is a preset threshold for the maximum tolerated degree of class imbalance ratio):

(a) Calculate the number of synthetic data examples that need to be generated for the minority class:

$$G = (m_l - m_s) \times \beta \quad (2)$$

Where $\beta \in [0, 1]$ is a parameter used to specify the desired balance level after generation of the synthetic data. $\beta = 1$ means a fully balanced data set is created after the generalization process.

(b) For each example $x_i \in \text{minorityclass}$, find K nearest neighbors based on the Euclidean distance in n dimensional space, and calculate the ratio r_i defined as:

$$r_i = \Delta_i/K, \quad i = 1, \dots, m_s \quad (3)$$

where Δ_i is the number of examples in the K nearest neighbors of x_i that belong to the majority class, therefore $r_i \in [0, 1]$;

(c) Normalize r_i according to $\hat{r}_i = r_i / \sum_{i=1}^{m_s} r_i$, so that \hat{r}_i is

a density distribution ($\sum_i \hat{r}_i = 1$)

(d) Calculate the number of synthetic data examples that need to be generated for each minority example x_i :

$$g_i = \hat{r}_i \times G \quad (4)$$

where G is the total number of synthetic data examples that need to be generated for the minority class as defined in Equation (2).

(e) For each minority class data example x_i , generate g_i synthetic data examples according to the following steps:

Do the **Loop** from 1 to g_i :

(i) Randomly choose one minority data example, x_{zi} , from the K nearest neighbors for data x_i .

(ii) Generate the synthetic data example:

$$s_i = x_i + (x_{zi} - x_i) \times \lambda \quad (5)$$

where $(x_{zi} - x_i)$ is the difference vector in n dimensional spaces, and λ is a random number: $\lambda \in [0, 1]$.

End **Loop**

Рисунок 1.4-Псевдокод методу ADASYN(джерело

[\[https://www.researchgate.net/figure/ADASYN-Pseudocode-16_fig2_367618805\]](https://www.researchgate.net/figure/ADASYN-Pseudocode-16_fig2_367618805))

Алгоритм ADASYN по суті є розширенням методу SMOTE. Він працює подібно до SMOTE, але додатково враховує густоту розподілу прикладів. Для цього алгоритм використовує k -найближчих сусідів та збільшує кількість прикладів шляхом додавання випадкового коефіцієнту до відстані між вибраним прикладом та його найближчим сусідом. Це дозволяє генерувати більше синтетичних прикладів для меншої класу, які знаходяться в менш щільно заміненіх областях, а менше для тих, які знаходяться в більш щільно заміненіх областях. Це покращує ефективність балансування класів та зменшує кількість зайвих синтетичних прикладів.

Незважаючи на те, що ADASYN є одним з найбільш популярних методів боротьби з незбалансованими класами, але є кілька недоліків:

1. Чутливість до шуму: метод ADASYN генерує нові приклади шляхом зміщення інших прикладів уздовж ліній від одного прикладу до іншого. Це може призвести до того, що генеровані приклади будуть більш схожі на шумові або випадкові відхилення від оригінальних прикладів, що може погіршити якість класифікації.

2. Обчислювальна складність: метод ADASYN вимагає обчислення відстаней між кожним з міноритарних прикладів та їхніми K найближчими сусідами, що може бути обчислювально складно для великих даних.

3. Залежність від K : якщо K занадто мале, то метод може бути недостатньо чутливим до структури даних, а якщо K занадто велике, то може виникнути проблема з вибором найближчих сусідів для прикладів в густонаселених областях.

4. Залежність від порогу: якщо поріг максимального ступеня дисбалансу встановлено неправильно, то метод може генерувати занадто багато або занадто мало штучних прикладів, що може призвести до перенавчання або недонавчання моделі.

5. Відсутність гарантії вдосконалення результатів: метод ADASYN не гарантує покращення результатів класифікації та може привести до погіршення результатів у деяких випадках.

Ще одним методом балансування є Class weighting:

Class weighting (вагова налаштування класів)-метод зменшення впливу нерівномірної розподіленості класів в навчальному наборі на результати моделі. Коли класи нерівномірно розподілені, модель може навчитися краще прогнозувати домінуючий клас, залишаючи меншість класів неохопленою. В такому випадку важко визначити відносну важливість кожного класу, іноді невдалим рішенням є використання простої точності (accuracy) як метрики оцінки результатів.

Метод вагової настройки полягає в наданні більшої ваги меншому класу в процесі навчання. Для цього вводиться параметр ваги (class weight), який може бути присвоєний кожному класу окремо. Параметри ваги можуть бути обчислені автоматично на основі розподілення класів у навчальному наборі або задані вручну.

Під час навчання, класи з вищою вагою будуть мати більший вплив на процес оптимізації моделі, що дозволяє досягти кращих результатів для менших класів. Зокрема, вагова настройка може допомогти уникнути перенавчання (overfitting) на домінуючих класах і покращити точність прогнозування менших класів.

Недоліком методу вагової настройки є те, що він може призвести до зниження точності для домінуючого класу, якщо вага для меншого класу встановлена занадто високо. Також, при використанні вагової настройки необхідно бути обережним і враховувати, що вага класу може впливати на рішення моделі, тому її необхідно встановлювати з урахуванням конкретних вимог до результатів.

1.5 Методи групування даних

Методи групування даних, такі як кластеризація, відіграють важливу роль у різних аспектах аналізу даних і машинного навчання. Основні цілі та застосування методів групування даних включають:

1. Розуміння структури даних: Методи групування даних дозволяють виявити приховану структуру та патерни в наборі даних. Вони допомагають ідентифікувати схожі групи об'єктів або спостережень, що можуть мати подібні характеристики або взаємозв'язки.

2. Візуалізація даних: Кластеризація дозволяє візуалізувати великий набір даних, розмістивши його у вигляді груп або кластерів. Це допомагає сприйняти та зрозуміти структуру даних, а також виявити аномалії або випадки, які виходять за межі груп.

3. Стиснення даних: Кластеризація може використовуватись для стиснення даних, замінюючи групи даних на їх центри або представників. Це може зменшити розмір даних та скоротити обчислювальні витрати при подальшому аналізі.

4. Класифікація та передбачення: Методи групування даних можуть використовуватись як попередній етап для побудови моделей класифікації або передбачення. Вони можуть допомогти виділити групи даних з подібними властивостями, що сприяє покращенню якості моделей.

6. Виявлення аномалій та шахрайства: Методи групування даних можуть виявляти аномалії або викиди в наборі даних, включаючи випадки шахрайства або несправедливої діяльності. Вони допомагають виділити незвичайні зразки, що відрізняються від нормального поведінки.

Тепер розглянемо деякі методи кластеризації більш детально:

Кластеризація k -середніх (k -means clustering)-один з найпоширеніших алгоритмів групування даних. Вона використовується для розділення набору даних на кілька груп або кластерів на основі їх схожості за деякими ознаками. Основна ідея полягає в тому, щоб знайти центри кластерів (k -середні) та призначити кожен об'єкт до найближчого за відстанню центру.

Основні кроки алгоритму k -середніх такі:

1. Вибрати кількість кластерів (k), яку потрібно сформувати з набору даних.
2. Випадковим чином ініціалізувати початкові k -середні (центри кластерів).

3. Повторювати до збіжності або до досягнення максимальної кількості ітерацій:

- Призначити кожен об'єкт до найближчого за відстанню центру кластеру.
- Оновити центри кластерів, обчислюючи середні значення об'єктів у кожному кластері.

4. Вивести остаточний набір кластерів та їх центри.

Кластеризація k -середніх має декілька переваг, зокрема простоту в реалізації, швидкість обчислень і підходить для великих наборів даних. Вона може використовуватись для виявлення прихованих структур та групування схожих об'єктів. Проте, важливо враховувати, що результати кластеризації можуть залежати від початкових значень центрів кластерів і можуть бути чутливі до шуму та випадкових факторів.

Кластеризація k -середніх може бути застосована в різних областях, таких як аналіз даних, комп'ютерний зір, біоінформатика, маркетингові дослідження, виявлення шахрайств та багато інших, де важливо групувати подібні об'єкти для подальшого аналізу та виявлення закономірностей.

Нижче наданий псевдокод алгоритму k -середніх зі статті [5]:

A. *K means Algorithm*

```

Input:
    D = {d1, d2, ..., dn} //set of n data items.
    k // Number of desired cluster

Output:
    A set of k clusters.

Step:
    1. Arbitrarily choose k data-items from D as
       initial centroids;
    2. Repeat

       Assign each item di to the cluster which
       has the closest centroid;
       Calculate new mean for each cluster;

    Until convergence criteria is met.
  
```

Fig. 1 Pseudocode of k means algorithm

Рисунок 1.5-Псевдо код алгоритму k-середніх

Нечітка кластеризація (fuzzy clustering) є іншим методом групування даних, який використовується для призначення кожного об'єкта набору даних до кластеру з певною ймовірністю членства. У відміню від кластеризації k-середніх, де об'єкт належить лише до одного кластеру з абсолютною впевненістю, в нечіткій кластеризації об'єкти можуть мати градацію впевненості у своєму членстві до різних кластерів.

Основні особливості нечіткої кластеризації:

1. Визначення кількості кластерів (k), аналогічно до кластеризації k-середніх.
2. Визначення функції належності, яка вказує ступінь членства кожного об'єкта до кожного кластеру.
3. Побудова матриці належності, де кожен елемент вказує ступінь членства об'єкта відносно кожного кластеру.
4. Оновлення центрів кластерів з урахуванням ступенів належності об'єктів.
5. Повторення процесу до збіжності або до заданої умови зупинки.

Нечітка кластеризація дозволяє моделювати неоднозначність та невизначеність при групуванні даних, що корисно в задачах, де об'єкти можуть мати часткову приналежність до кількох кластерів одночасно. Цей метод знайшов застосування в багатьох областях, включаючи розпізнавання образів, обробку природних мов, біоінформатику та інші. Важливо враховувати, що нечітка кластеризація вимагає визначення додаткових параметрів, таких як функція належності, і може бути чутливою до початкових значень.

Нижче наданий псевдокод алгоритму нечіткої кластеризації зі статті [6]:

```

FCM Algorithm :
Input :  $\theta (k,1), N$ 
Output  $U^*_{FCM}, V^*_{FCM}$ 
1. Initialize Partition  $U^{(0)}$  randomly
2. for  $i = 1$  to  $n$ 
3.   for  $k = 1$  to  $c$ 
4.     Repeat for  $j = 1, 2, 3, \dots$ 
5.       Update centroid  $V^{(j)}$  with  $U^{(j-1)}$  Using (3)
6.       Compute Distance  $D^{(j)}$  with  $V^{(j)}$ 
7.       Update Partition Matrix  $U^{(j)}$  with  $D^{(j)}$  using (5)
8.       Until  $\|U^{(j)} - U^{(j-1)}\| < \epsilon$ 
9.     end
10.  end
11. Return  $U^*_{FCM} \leftarrow U^{(j)}$  and  $V^*_{FCM} \leftarrow V^{(j)}$ 

```

Рисунок 1.6-Псевдокод алгоритму нечіткої кластеризації

1.6 Методи класифікації незбалансованих наборів даних

XGBoost (eXtreme Gradient Boosting)-алгоритм машинного навчання з відкритим кодом, який використовує градієнтний бустінг для підвищення точності моделей класифікації та регресії. Він є популярним в індустрії та наукових дослідженнях завдяки своїй ефективності та швидкості.

XGBoost є розширенням градієнтного бустінгу, додавши багато нових функцій та покращень, таких як розріджена матриця, підтримка стандартних API для моделей машинного навчання та багато іншого. В основі XGBoost лежить бібліотека на мові програмування C++, яка підтримує паралельну обробку даних на рівні батчів.

Алгоритм XGBoost використовує градієнтний бустінг над деревами рішень. Основна ідея полягає в тому, щоб послідовно побудувати послідовні дерева рішень і використовувати їх для корекції попередніх помилок. Кожне наступне дерево побудоване на основі остач, залишених попереднім деревом.

Алгоритм має декілька параметрів, які можна настроїти для отримання найкращих результатів. Наприклад, кількість дерев, що будуються, глибина дерев, критерії розбиття, розмір батчу та інші.

Основна ідея полягає в тому, щоб послідовно додавати нові дерева до ансамблю, коригуючи попередні прогнози моделі, щоб отримати кращу точність прогнозів.

Основний алгоритм XGBoost можна розбити на наступні кроки:

1. Ініціалізувати модель: визначити стартове прогнозоване значення для всіх прикладів і створити модель, яка буде постійно оновлюватись.

2. Обчислити градієнт та гессіан для всіх прикладів: градієнт і гессіан-це вектори, які вказують на напрямок найшвидшого зростання функції в кожній точці. Градієнт використовується для покращення моделі на кожній ітерації, а гессіан допомагає підвищити ефективність покращення моделі.

3. Побудувати дерево рішень: дерево будується шляхом розділення даних на підмножини, використовуючи критерій інформативності (наприклад, критерій Джині або ентропія). Алгоритм рекурсивно розділяє дані на підмножини, поки не досягне максимальної глибини дерева або не досягне мінімальної кількості елементів в кожному листку.

4. Розрахувати приналежність до кожного листка дерева для кожного прикладу: кожен приклад розподіляється на шляху вниз по дереву, завдяки чому кожен елемент даних потрапляє в листок, до якого належить.

5. Обчислити розбіжність між прогнозованими та фактичними значеннями: різниця між прогнозованим значенням та фактичним значенням цільової змінної для кожного елемента даних.

6. Оновити модель, додаючи нове дерево до ансамблю: нове дерево додається до моделі, і прогнозовані значення оновлюються з урахуванням нової інформації. Цей процес повторюється доти, доки не досягнута задана кількість дерев або поки досягнута задана точність.

7. Здійснити прогноз: для нових елементів даних здійснюється прогноз за допомогою побудованого ансамблю дерев.

Основними перевагами XGBoost є швидкість та ефективність роботи з великими наборами даних, а також можливість автоматичного вибору гіперпараметрів, що дозволяє покращити точність моделі. Крім того, XGBoost має

вбудовану підтримку обробки пропущених даних та можливість роботи з різними типами даних.

Недоліками XGBoost можуть бути:

- складність налаштування параметрів моделі;
- необхідність розуміння технічних аспектів реалізації.

Нижче представлено псевдокод зі статті ADASYN: Adaptive Synthetic Sampling Approach for Imbalanced Learning [2]

- Initialize sample weights $w_i^{(0)} = \frac{1}{l}, i = 1, \dots, l$.
- For all $t = 1, \dots, T$
 - Train base algo b_t , let ϵ_t be it's training error.
 - $\alpha_t = \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$.
 - Update sample weights: $w_i^{(t)} = w_i^{(t-1)} e^{-\alpha_t y_i b_t(x_i)}, i = 1, \dots, l$.
 - Normalize sample weights: $w_0^{(t)} = \sum_{j=1}^k w_j^{(t)}, w_i^{(t)} = \frac{w_i^{(t)}}{w_0^{(t)}}, i = 1, \dots, l$.
- Return $\sum_t \alpha_t b_t$

Рисунок 1.7-Псевдокод XGBoost

Узагальнюючи, XGBoost є потужним та ефективним алгоритмом машинного навчання, який дозволяє вирішувати різноманітні задачі з високою точністю. Однак, перед використанням алгоритму, необхідно ретельно налаштувати його параметри та врахувати його технічні особливості.

CatBoost (Categorical Boosting)-алгоритм машинного навчання, що використовує градієнтний бустинг для задач класифікації та регресії, розроблений для обробки як числових, так і категоріальних ознак без їхньої попередньої обробки. Завдяки своїй продуктивності та ефективності, він широко використовується в індустрії та наукових дослідженнях, особливо для задач, де є великий обсяг категоріальних даних.

CatBoost базується на градієнтному бустингу дерев рішень і має кілька особливостей, які відрізняють його від інших алгоритмів бустингу, таких як

XGBoost і LightGBM. CatBoost поєднує в собі метод Ordered Boosting і унікальний спосіб обробки категоріальних ознак, що підвищує стабільність і точність моделі.

Основні етапи роботи CatBoost:

1. Ініціалізація моделі: визначається початковий прогноз для всіх прикладів, який буде поступово покращуватись з кожною ітерацією побудови дерев.

2. Обробка категоріальних ознак: на відміну від інших алгоритмів, CatBoost автоматично перетворює категоріальні ознаки на числові, застосовуючи Target Encoding (кодування на основі значення цільової змінної) або Leave-One-Out Encoding, що знижує ризик переобучення.

3. Ordered Boosting: унікальна технологія CatBoost, що запобігає переобученню, використовуючи "упорядкований" бустинг. На кожній ітерації при побудові дерев використовуються нові підмножини даних, завдяки чому кожне дерево будується на основі незалежних значень, що покращує узгодженість моделі.

4. Побудова дерева рішень: дерева рішень створюються за принципом рекурсивного поділу даних на підмножини, щоб мінімізувати функцію втрат. Критерії поділу, такі як інформаційна вигода або ентропія, допомагають обрати оптимальні точки поділу.

5. Обчислення прогнозів для кожного листка: кожен приклад у наборі даних проходить по шляху дерева, потрапляючи у відповідний листок, де отримує певне прогнозне значення.

6. Оновлення ансамблю: побудоване дерево додається до поточної моделі, і прогнозовані значення оновлюються на основі залишкових помилок, які визначаються як різниця між фактичними та прогнозованими значеннями.

7. Циклічне повторення: модель продовжує додавати нові дерева до ансамблю, доки не досягне заданої кількості ітерацій або не досягне потрібного значення точності.

CatBoost має різноманітні параметри, які можна налаштувати для досягнення оптимальної продуктивності:

- `iterations`: кількість ітерацій навчання або кількість дерев, які буде додано до ансамблю.
- `learning_rate`: впливає на швидкість навчання моделі, менш інтенсивний темп може забезпечити кращу узгодженість за рахунок більшої кількості ітерацій.
- `depth`: максимальна глибина кожного дерева, що впливає на складність та точність моделі.
- `l2_leaf_reg`: параметр регуляризації для уникнення переобучення.
- `bagging_temperature`: контролює випадковість під час вибору підмножин, що може поліпшити узагальнення, але зменшити стабільність.
- `border_count`: кількість розподілів для числових ознак, що впливає на обчислення значень.

Переваги CatBoost:

- Ефективна обробка категоріальних ознак: економить час на підготовці даних.
- Устойчивість до переобучення завдяки Ordered Boosting.
- Висока продуктивність і паралельність обробки: алгоритм оптимізований для багатоядерних процесорів.
- Мала чутливість до налаштування гіперпараметрів: базові параметри часто забезпечують хороші результати, особливо для задач з великим числом ознак.

Недоліки CatBoost:

- Вимоги до обчислювальних ресурсів: на дуже великих наборах даних може вимагати значної оперативної пам'яті.
- Складність інтерпретації результатів: через велику кількість побудованих дерев та відсутність чітких правил інтерпретації результатів.

Нижче приведено псевдокод алгоритму CatBoost зі статті «A NEW APPROACH TO CUSTOMER CHURN MODELLING USING LLM AND CATBOOST» [3]:

Algorithm 1 CatBoost Training Pseudocode

Require: Training data: $X_{\text{train}}, y_{\text{train}}$
Require: Number of iterations (trees): M
Require: Learning rate: η
Require: Number of tree nodes: K
Require: Categorical features: cat_features

```

0: Initialize:  $F(x) = 0$ 
0: for  $m = 1$  to  $M$  do
0:   Compute negative gradient for each sample:
0:   for  $i = 1$  to  $N$  do
0:     Compute:  $g_i = -\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)}$ 
0:   end for
0:   Train a decision tree on the negative gradients:
0:   Use  $\text{cat\_features}$  for handling categorical variables.
0:   Tree structure is determined by minimizing the loss
    function.
0:   Each leaf node corresponds to a prediction value for a
    subset of samples.
0:   Compute optimal weight for the new tree:
0:   Compute:  $\gamma_m = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i))$ 
0:   Update the model:
0:   Update:  $F(x) = F(x) + \eta \gamma_m h_m(x)$ 
0: end for

```

Рисунок 1.9-Псевдокод CatBoost

CatBoost є потужним інструментом для роботи з категоріальними даними та застосовується в різних галузях. Завдяки унікальній обробці категоріальних ознак та Ordered Boosting він демонструє високу точність та продуктивність на багатьох реальних наборах даних, проте для досягнення максимальних результатів потрібне правильне налаштування параметрів.

LightGBM (Light Gradient Boosting Machine)—бібліотека машинного навчання, що використовує градієнтний бустинг на базі дерев рішень. Вона була розроблена компанією Microsoft і відома своєю швидкістю та ефективністю, особливо для великих і високовимірних наборів даних. LightGBM часто використовується в задачах класифікації, регресії та ранжування завдяки своїй здатності швидко обробляти великі обсяги даних і забезпечувати високу точність.

Основною відмінністю LightGBM є те, що він використовує інноваційний підхід під назвою Leaf-Wise Tree Growth (зростання дерева по листкам), на відміну від традиційного підходу рівневого зростання. Завдяки цьому LightGBM може краще узагальнювати дані, що підвищує точність моделі при меншій глибині дерева і обмежує переобучення.

Основні етапи роботи LightGBM:

1. Ініціалізація моделі: задається початковий прогноз для всіх прикладів. На кожній ітерації модель додає нове дерево, щоб зменшити залишкові помилки.
2. Leaf-Wise Growth (зростання дерева по листкам): кожне нове дерево в LightGBM будується шляхом додавання нових листків до існуючої структури. Це означає, що замість того, щоб збалансовано розширювати кожен рівень дерева, LightGBM вибирає і розширює найбільш інформаційно насичені гілки. Цей підхід дозволяє моделі навчатися швидше, але іноді може призводити до переобучення.
3. Histogram-Based Splitting (поділ на основі гістограм): LightGBM використовує метод розподілу значень у вигляді гістограм для кожного поділу. Це значно знижує обчислювальні витрати та дозволяє швидше визначити оптимальні точки поділу, навіть на великих наборах даних.
4. Data Partitioning and Parallel Learning (розподіл даних та паралельне навчання): LightGBM підтримує паралельне навчання, розподіляючи дані на підмножини і виконуючи обчислення паралельно. Це дозволяє ефективніше використовувати обчислювальні ресурси.

5. Обчислення градієнтів та оновлення ансамблю: на кожній ітерації обчислюються градієнти та залишкові помилки для кожного прикладу. Нове дерево додається до ансамблю, оновлюючи прогнозовані значення для кожного елемента даних.
6. Циклічне повторення: процес триває доти, доки не досягнуто заданої кількості ітерацій або не досягнуто необхідної точності.

Основні параметри налаштування LightGBM

- `num_leaves`: кількість листків у кожному дереві. Велика кількість листків підвищує складність і точність моделі, але може спричинити переобучення.
- `learning_rate`: темп навчання моделі. Низьке значення `learning_rate` допомагає зменшити похибки за рахунок більшої кількості ітерацій.
- `max_depth`: максимальна глибина дерева. Встановлюється для запобігання переобученню.
- `min_data_in_leaf`: мінімальна кількість прикладів у листку. Параметр допомагає обмежити розмір листка, що сприяє узагальненню.
- `feature_fraction`: частка ознак, що використовуються при побудові кожного дерева. Допомагає знизити ризик переобучення.
- `bagging_fraction` та `bagging_freq`: встановлюють частку даних для навчання на кожній ітерації та частоту її вибірки.

Переваги LightGBM:

- Висока швидкість навчання та прогнозування: LightGBM швидше обробляє великі набори даних завдяки Leaf-Wise Growth та Histogram-Based Splitting.
- Ефективність на великих наборах даних: LightGBM добре підходить для обробки великих і складних датасетів із численними ознаками.
- Мала чутливість до пропусків в даних: LightGBM добре працює навіть із відсутніми значеннями.

Недоліки LightGBM:

- Схильність до переобучення: Leaf-Wise Growth, хоч і дозволяє досягти високої точності, може призводити до переобучення, особливо якщо кількість листків (`num_leaves`) надмірно велика.

- Складне налаштування параметрів: потребує ретельного налаштування для уникнення переобучення та забезпечення оптимальної продуктивності.

Нижче представлено псевдокод алгоритму зі статті «A Novel Prediction Model for Diabetes Detection Using Gridsearch and A Voting Classifier between Lightgbm and KNN»[4]:

The LightGBM algorithm

Input:

Training data: $D = \{(\chi_1, y_1), (\chi_2, y_2), \dots, (\chi_N, y_N)\}$, $\chi_i \in \mathcal{X}$, $\mathcal{X} \subseteq \mathbb{R}$, $y_i \in \{-1, +1\}$; loss function: $L(y, \theta(\chi))$;

Iterations:

M ; Big gradient data sampling ratio: a ; slight gradient data sampling ratio: b ;

1: Combine features that are mutually exclusive (i.e., features never simultaneously accept nonzero values) of χ_i , $i = \{1, \dots, N\}$ by the exclusive feature bundling (EFB) technique;

2: Set $\theta_0(\chi) = \arg \min_c \sum_i^N L(y_i, c)$;

3: For $m = 1$ to M do

4: Calculate gradient absolute values:

$$r_i = \left| \frac{\partial L(y_i, \theta(x_i))}{\partial \theta(x_i)} \right|_{\theta(x) = \theta_{m-1}(x)}, i = \{1, \dots, N\}$$

5: Resample data set using gradient-based one-side sampling (GOSS) process:

$topN = a \times len(D)$; $randN = b \times len(D)$;

$sorted = GetSortedIndices(abs(r))$;

$A = sorted[1 : topN]$; $B = RandomPick(sorted[topN : len(D)], randN)$;

$D' = A + B$;

6: Calculate information gains:

$$V_j(d) = \frac{1}{n} \left(\frac{\left(\sum_{x_i \in A_l} r_i + \frac{1-a}{b} \sum_{x_i \in B_l} r_i \right)^2}{n_l^j(d)} + \frac{\left(\sum_{x_i \in A_r} r_i + \frac{1-a}{b} \sum_{x_i \in B_r} r_i \right)^2}{n_r^j(d)} \right)$$

7: Develop a new decision tree $\theta_m(x)'$ on set D'

8: Update $\theta_m(\chi) = \theta_{m-1}(\chi) + \theta_m(\chi)$

9: End for

10: Return $\hat{\theta}(x) = \theta_M(x)$

Рисунок 1.9-Псевдокод LightGBM

LightGBM є потужним інструментом для задач класифікації, регресії та ранжування, що дозволяє ефективно обробляти великі та складні набори даних з численними ознаками. Завдяки своїй швидкості та можливості обробки даних із відсутніми значеннями, він широко використовується в задачах, де необхідна висока продуктивність і точність.

1.7 Методи оптимізації параметрів моделей

Оптимізація гіперпараметрів для моделей XGBoost, LightGBM та CatBoost є критичною для досягнення максимальної ефективності в задачах машинного навчання, особливо в умовах великих даних або незбалансованих класів. У таких завданнях, як виявлення шахрайства, де важливі навіть незначні помилки класифікації, правильно налаштовані параметри можуть суттєво підвищити якість передбачень, наприклад, зменшивши частку помилкових спрацювань або покращивши розпізнавання рідкісних випадків. Оскільки кожна з моделей (XGBoost, LightGBM, CatBoost) має десятки параметрів, які впливають на її поведінку та ефективність, важливо знайти оптимальні комбінації для конкретної задачі. Це вимагає спеціальних методів оптимізації, таких як Optuna, Random Search, Grid Search та інші підходи. Вони значно полегшують процес пошуку найкращих параметрів і дозволяють зосередитися на досягненні ключових показників продуктивності, таких як F1-скор, точність або AUC.

Основні методи оптимізації гіперпараметрів

1. Optuna:

Optuna є однією з найпотужніших бібліотек для автоматизованої оптимізації гіперпараметрів, заснованої на байєсівському підході та техніці під назвою Tree-structured Parzen Estimator (TPE). Вона здійснює пошук параметрів в інтерактивному режимі, що дозволяє змінювати діапазони параметрів "на ходу" та підвищувати гнучкість процесу оптимізації. Основні переваги Optuna:

- Підтримка байєсівської оптимізації, що скорочує кількість необхідних ітерацій;
- Можливість паралельного пошуку, що особливо корисно для великих наборів даних;
- Підтримка розумного вибору параметрів для різних моделей, таких як XGBoost, LightGBM та CatBoost;
- Динамічна зміна простору пошуку, що дозволяє адаптувати оптимізацію під конкретну модель або метрику.

2. Random Search:

Random Search передбачає випадковий вибір значень параметрів з визначених діапазонів. Це простий і швидкий метод, який, незважаючи на свою випадковість, є ефективним на початкових етапах оптимізації. У випадках, коли простір параметрів великий або часу на вичерпний пошук недостатньо, Random Search може допомогти швидко виявити загальні тенденції. Він також добре масштабується і дозволяє швидко протестувати різні комбінації параметрів. Однак цей метод не завжди приводить до оптимального рішення, особливо коли параметри сильно залежать один від одного.

3. Grid Search:

Grid Search є систематичним методом, який перебирає всі можливі комбінації параметрів у визначеному діапазоні. Він гарантує, що буде знайдено найкращий набір параметрів з усіх можливих варіантів, однак має значний недолік — його обчислювальні витрати зростають експоненційно з кількістю параметрів. Grid Search зазвичай використовується тоді, коли:

- Потрібно знайти точний локальний оптимум;
- Кількість параметрів обмежена, а діапазони параметрів невеликі;
- Час обчислень та ресурси дозволяють здійснити вичерпний перебір.

4. Байєсівська оптимізація (Bayesian Optimization):

Це складніший підхід, який передбачає побудову моделі функції оптимізації на основі попередніх результатів і використання цієї моделі для вибору наступних наборів параметрів. Байєсівська оптимізація дозволяє ефективно зосередитися на діапазонах параметрів, які мають найвищий потенціал покращення, що допомагає скоротити час пошуку оптимальних параметрів. Цей метод корисний у задачах, де кожна ітерація потребує значного обчислювального ресурсу.

5. Nuropt та його TPE-алгоритм:

Nuropt, подібно до Optuna, реалізує Tree-structured Parzen Estimator (TPE) для байєсівської оптимізації. Він є потужним інструментом, який допомагає швидко здійснювати пошук у великому просторі параметрів. Цей підхід корисний для

моделей, які потребують великої кількості налаштувань, і може використовувати паралельні обчислення, що значно скорочує час оптимізації.

6. Evolutionary Algorithms:

Еволюційні алгоритми, такі як Genetic Algorithm, використовують принципи природного відбору, щоб знайти найкращі параметри. Вони імітують процес еволюції, створюючи нові набори параметрів на основі найкращих комбінацій з попередніх ітерацій. Такий підхід добре працює з дуже великими або складними просторами параметрів, де інші методи виявляються надто повільними або неефективними.

Кожен з цих методів має свої особливості, і правильний вибір підходу до оптимізації залежить від конкретної задачі, обсягу даних і доступних обчислювальних ресурсів. Використання таких алгоритмів для пошуку гіперпараметрів дозволяє значно поліпшити результати моделей XGBoost, LightGBM і CatBoost, роблячи їх більш ефективними, точними й надійними в практичних завданнях.

1.8 Показники ефективності моделей класифікації

Показники ефективності моделей класифікації визначають, наскільки точно модель може передбачити клас об'єкта. Існує багато різних показників, і кожен з них дає певну інформацію про ефективність моделі.

Матриця помилок (Confusion Matrix)-таблиця, яка демонструє кількість правильних та неправильних прогнозів, зроблених моделлю класифікації на основі тестових даних. Вона дозволяє зрозуміти, наскільки добре модель працює на різних класах.

Матриця помилок містить чотири елементи:

- True Positive (TP) - кількість правильно визначених позитивних випадків;
- False Positive (FP) - кількість неправильно визначених позитивних випадків;
- True Negative (TN) - кількість правильно визначених негативних випадків;

- False Negative (FN) - кількість неправильно визначених негативних випадків.

Нижче надано двокласову матрицю у вигляді таблиці:

		Predicted	
		Negative (N) -	Positive (P) +
Actual	Negative -	True Negative (TN)	False Positive (FP) Type I Error
	Positive +	False Negative (FN) Type II Error	True Positive (TP)

Рисунок 1.10-Вигляд двокласової confusion matrix у вигляді таблиці

Та вигляд багатокласової матриці:

		PREDICTED classification			
		Classes	a	b	c
ACTUAL classification	a	TN	FP	TN	TN
	b	FN	TP	FN	FN
	c	TN	FP	TN	TN
	d	TN	FP	TN	TN

Рисунок 1.11-Вигляд багатокласової confusion matrix у вигляді таблиці

За допомогою матриці помилок можна розрахувати різні метрики ефективності моделі, такі як точність, чутливість, специфічність та F1-оцінка.

Нижче надані розвідкова інформація щодо цих метрик::

1. Accuracy - відношення кількості правильно класифікованих об'єктів до загальної кількості об'єктів у вибірці.

Формула точності:

$$accuracy = (TP + TN) / (TP + TN + FP + FN) \quad (1.1)$$

де TP - кількість правильно визначених позитивних випадків, TN - кількість правильно визначених негативних випадків, FP - кількість неправильно визначених позитивних випадків, а FN - кількість неправильно визначених негативних випадків.

Недоліком accuracy є те, що вона може бути недостатньо інформативною у випадках, коли класи незбалансовані, тобто коли один клас має значно меншу кількість екземплярів, ніж інший. У таких випадках accuracy може бути високою, навіть якщо модель погано працює на меншому класі. Наприклад, якщо ми маємо датасет з 1000 зразків, з яких 900 належать до класу А, а 100 - до класу В, то класифікатор може досить легко досягти accuracy 90% за рахунок того, що він буде правильно класифікувати зразки класу А, але може погано працювати на зразках класу В, які він буде помилково віднести до класу А. Таким чином, метрика accuracy не відображає реальної ефективності моделі у випадку з незбалансованими даними.

2. Precision-метрика, яка використовується для оцінки того, наскільки часто модель правильно класифікує екземпляри позитивного класу. Ця метрика розраховується як відношення кількості правильно класифікованих позитивних екземплярів до загальної кількості позитивних екземплярів, що були класифіковані моделлю:

$$Precision = TP / (TP + FP) \quad (1.2)$$

Precision часто використовується в задачах, де важливо не допустити помилкових позитивних прогнозів. Наприклад, в медицині, якщо модель передбачає наявність захворювання, але насправді воно відсутнє, то це може призвести до неправильного лікування. У такому випадку ми більше зацікавлені у високій precision, ніж у високій чутливості (recall).

Однак, precision може бути обманливою метрикою, особливо в задачах з незбалансованими даними, коли кількість екземплярів позитивного класу дуже мала порівняно з екземплярами негативного класу. У таких випадках модель може бути дуже точною в класифікації негативних екземплярів, але не здатна коректно класифікувати позитивні. Тому, перед використанням precision як метрики, необхідно детально проаналізувати дані та визначити, чи є вони збалансованими, і чи необхідно використовувати інші метрики для оцінки ефективності моделі.

3. Повнота (Recall), також відома як чутливість (sensitivity) або True Positive Rate (TPR), є метрикою, що вимірює здатність класифікатора правильно визначати позитивні зразки з усіх дійсно позитивних зразків. Вона вказує на те, яка частка позитивних зразків була виявлена класифікатором.

Повнота обчислюється за формулою:

$$Recall = TP / (TP + FN) \quad (1.3)$$

Це відношення показує, яка частка позитивних зразків була виявлена класифікатором. Чим більше значення повноти, тим краще класифікатор виявляє позитивні зразки. Високе значення повноти означає, що класифікатор мало пропускає позитивні зразки, але водночас може призводити до більшої кількості неправильно класифікованих негативних зразків.

Повнота є важливою метрикою в задачах, де виявлення позитивних зразків має високий пріоритет, наприклад, у виявленні хвороб, де не діагноз може мати серйозні наслідки.

4. F1-score (F1-міра) є гармонічним середнім між точністю та повнотою, що використовується для вимірювання точності бінарної класифікації. Це є зважено середнє між точністю і повнотою, і він приймає значення між 0 та 1.

Формула F1-score:

$$F1 - score = 2 * (precision * recall) / (precision + recall) \quad (1.4)$$

F1-міра використовується в тих випадках, коли точність та повнота мають різне значення. Якщо точність висока, а повнота низька, то значення F1-міри буде низьким, а якщо точність низька, а повнота висока, то F1-міра також буде низькою. F1-міра зазвичай використовується, коли дані є незбалансованими,

тобто один клас має значно більше прикладів, ніж інший. Наприклад, якщо ми маємо 1000 зображень, і 950 з них належать до класу "кіт", а 50 до класу "собака", то модель може класифікувати всі зображення як "кіт", отримуючи точність 95%. Однак, це буде неправильною моделлю, тому що вона не може відрізнити між котами та собаками. У цьому випадку важливішою буде повнота, оскільки ми хочемо, щоб модель змогла визначати обидва класи. F1-міра враховує як точність, так і повноту, тому вона є кращою метрикою в таких випадках.

Недоліком F1-міри є те, що вона не враховує зразки, які були класифіковані правильно, але помилково визначені як інший клас. В такому випадку метрика може бути низькою, навіть якщо модель може коректно визначати більшість зразків. Також F1-міра не дозволяє враховувати більшість поганих результатів в одному з класів, що може бути важливою інформацією у деяких задачах. Наприклад, якщо ми працюємо з даними про пацієнтів з деяким захворюванням, то помилкове визначення хворих пацієнтів як здорових може бути критично важливим.

Перевагою F1-міри є те, що вона є зваженою мірою, яка враховує обидві метрики - точність та повноту. Це дозволяє отримати більш точну інформацію про ефективність моделі в задачах класифікації з незбалансованими даними. Крім того, F1-міра є доброю метрикою, коли точність та повнота мають приблизно однакове значення.

Загалом, F1-міра є корисною метрикою, яка дозволяє враховувати як точність, так і повноту, та допомагає визначити ефективність моделі в задачах класифікації з незбалансованими даними.

6. AUC-ROC (Area Under the Receiver Operating Characteristic Curve)-метрика, яка використовується для оцінки якості бінарної моделі класифікації. Ця метрика вимірює площу під кривою ROC (Receiver Operating Characteristic), яка є графіком, що показує залежність між True Positive Rate (TPR) і False Positive Rate (FPR) при зміні порога відсічення.

TPR –вже відома нам метрика, Повнота (Recall) (1.3), яку ми вже розглянули. FPR-відношення неправильно класифікованих негативних прикладів до загальної кількості негативних прикладів в тестовому наборі.

Формули для AUC-ROC:

$$TRP = \frac{TP}{TP+FN}$$

(1.5)

$$FRP = FP/(FP + TN)$$

(1.6)

Крива ROC дає можливість оцінити те, як добре модель розділяє позитивні та негативні класи, в залежності від значення порога відсічення. Ідеальна модель буде мати площу під кривою дорівнює 1, тоді як модель, яка не відрізняється від випадкового вибору, матиме площу дорівнює 0.5.

$$AUC - ROC = \int_0^1 TPR(f(x))dFRP(fx)$$

(1.7)

AUC-ROC дуже корисна метрика в разі, коли класи в даних незбалансовані, тобто один з класів має більше прикладів, ніж інший. Оскільки ця метрика не залежить від пропорцій класів, вона може допомогти уникнути помилкового враження про те, що модель є досить точною, тоді як насправді вона дає низьку точність на менш представлених класах.

До переваг AUC-ROC можна віднести її стійкість до шуму та зміни порога відсічення, що дозволяє використовувати цю метрику в різних ситуаціях. Також важливою перевагою є те, що AUC-ROC можна використовувати для порівняння різних моделей, які можуть мати різні значення порога відсічення.

Недоліки AUC-ROC полягають у тому, що вона не дає змоги оцінити ефективність різних класифікаторів на основі різних порогів відсічення на прикладах з різною важливістю, оскільки площа під кривою є вагомою метрикою, і значення цієї метрики може бути неадекватним для визначення ефективності класифікатора на деяких конкретних важливих прикладах. Також AUC-ROC може

не підходити для задач, де дуже важливо, щоб модель правильно класифікувала певний клас, який є меншим за інші класи в тестовому наборі.

7. Logloss-метрика, яка використовується для оцінки точності класифікаторів. Вона є функцією відстані між прогнозованими й фактичними значеннями. Ця метрика вимірюється в діапазоні від 0 до нескінченності, де 0 позначає ідеальну точність, а значення, які прямують до нескінченності, вказують на повну невідповідність прогнозів реальності.

Формула для розрахунку log loss (1.6):

$$\log loss = -(1/N) * \sum [y * \log(p) + (1 - y) * \log(1 - p)] \quad (1.8)$$

де N - загальна кількість прикладів у наборі даних, y - фактична мітка класу (0 або 1), а p - передбачена ймовірність належності до класу 1 за моделлю.

Сума береться по всіх прикладах в наборі даних, і для кожного прикладу обчислюється вираз $y * \log(p) + (1 - y) * \log(1 - p)$. Потім отримані значення сумуються і помножуються на $-(1/N)$, що дає загальну логарифмічну втрату для набору даних.

Log loss доречно використовувати в задачах бінарної та багатокласової класифікації. Його головною перевагою є те, що він дуже чутливий до неправильних передбачень, тобто якщо класифікатор вважає один клас більш ймовірним, але насправді він належить до іншого класу, то це значно підвищить значення log loss. Завдяки цьому метрика дає важливу інформацію про те, як далеко класифікатор відхиляється від правильних відповідей. Однак він також може бути чутливим до перевірочних даних та до великої кількості класів.

Кожен з цих показників має свої переваги і недоліки, тому важливо використовувати декілька показників для оцінки ефективності моделі і приймати рішення залежності від контексту задачі. Наприклад, якщо ми важливіше уникнути помилок визначення позитивних об'єктів, то точність і точність класифікації позитивних об'єктів (precision) можуть бути більш інформативними показниками. У той же час, якщо ми важливіше забезпечити, щоб не було

пропущених позитивних об'єктів, тоді повнота (recall) може бути більш інформативною метрикою.

Висновок

В ході аналізу та вивчення теоретичного матеріалу кваліфікаційної роботи були розглянуті наступні теми:

Аналіз та дослідження банківських транзакції та особливості їх набору даних, що впливають на методи та якість класифікації.

Методи підвищення якості прогнозу моделей класифікації, їх вплив на моделі, дані та результат, їх недоліки, переваги та чим кожен з них відрізняються від інших.

Також розглянуті були нерівномірність розподілу класів та проблема нерівномірної класифікації у задачах класифікації, їх вплив на моделі прогнозування.

Вивчені та проаналізовані методи балансування даних у незбалансованому наборі даних, такі як oversampling, undersampling, SMOTE та ADASYN, їх переваги, недоліки та в яких задачах їх краще використовувати.

Методи групування даних, такі як алгоритм K-середніх та нечіткої кластеризації, їх зміст, призначення до роботи та чим кожен з них відрізняються від іншого

Методи класифікації незбалансованих даних, такі як XGBoost та Random Forest.

Показники ефективності або метрики моделей класифікації такі як Accuracy, Precision, AUC-ROC, Log-loss та F1-score, проаналізовано та вивчено чим вони відрізняються та в яких задачах їх більш коректно застосовувати.

Після вивчення теоретичного матеріалу ми ознайомилися з проблемами незбалансованих даних у задачах класифікацій. Отриманий теоретичний базис надає можливість розробити ефективний підхід до вирішення задачі прогнозування шахрайства на основі незбалансованих даних. У наступних

розділах роботи буде розглянуто практичне застосування цих методів та технік для побудови класифікаційної моделі, яка зможе ефективно прогнозувати випадки шахрайських транзакцій враховуючи незбалансованість даних.

2 СПЕЦІАЛЬНИЙ РОЗДІЛ

2.1 Постановка задачі

Задача полягає у розробці та порівняння моделей машинного навчання для виявлення мобільних шахрайських транзакцій на основі синтетичного незбалансованого набору даних "transactions_train.csv".

2.1.1 Мета дослідження

Метою даного дослідження полягає в застосуванні та порівнянні методів класифікації незбалансованих наборів даних для використання в задачах виявлення фроду. Конкретніше, метою є аналіз та порівняння методів XGBoost з двома альтернативними методами прогнозування транзакцій на наявність фроду на основі незбалансованих даних, таких як CatBoost та LightGBM. Дослідження має мету зрозуміти, чи можуть ці методи класифікації бути ефективнішими при роботі з незбалансованими даними в контексті задач виявлення фроду, ніж XGBoost.

2.1.2 Задачі дослідження

Для досягнення поставленої мети необхідно вирішити такі задачі:

1. Розглянути теоретичні аспекти класифікації даних, зокрема, методи роботи з незбалансованими даними.
2. Оглянути різноманітні методи класифікації незбалансованих наборів даних, включаючи методи, які засновані на деревах рішень, нейронних мережах, ансамблевих методи та інші.

3. Дослідити існуючі наукові роботи, присвячені використанню методів класифікації для задач виявлення фрода.

4. Описати методи XGBoost, CatBoost, LightGBM та порівняти їх.

5. Розробити експериментальну частину, в рамках якої буде проведено експеримент з використанням методу XGBoost, CatBoost та LightGBM на даних з задачі виявлення фрода.

6. Проаналізувати результати експерименту та зробити висновки про ефективність методів XGBoost, CatBoost та LightGBM у класифікації незбалансованих наборів даних в задачах виявлення фрода.

7. Зробити висновки щодо зробленої роботи та актуальності та практичних застосувань цих методів у задачах на виявлення фроду.

2.2 Обґрунтування вибору середовища програмування

Для прогнозування фроду (шахрайства) існує багато середовищ програмування, які можна використовувати. Вибір конкретного середовища залежить від ваших вимог, знань та вподобань. Ось декілька популярних середовищ програмування, які часто використовуються для прогнозування фроду:

1. Python: Python є однією з найпоширеніших мов програмування для аналізу даних та машинного навчання. Він має багато потужних бібліотек, таких як scikit-learn, TensorFlow, PyTorch і XGBoost, які допомагають у побудові моделей прогнозування фроду.

2. R - мова програмування та середовище для статистичного аналізу та візуалізації даних. Вона також має багато пакетів для машинного навчання, таких як caret, randomForest і xgboost, які можуть бути використані для прогнозування фроду.

3. MATLAB - інтерактивна система для чисельних обчислень та аналізу даних. Вона має багато інструментів для машинного навчання, таких як Statistics and Machine Learning Toolbox, які можуть бути використані для прогнозування фроду.

4. SAS - програмний пакет для статистичного аналізу та аналізу даних. Він має спеціалізовані інструменти для прогнозування шахрайства, такі як SAS Fraud Framework, який надає засоби для виявлення та прогнозування шахрайства.

5. Spark - це фреймворк для обробки великих обсягів даних та аналізу в розподіленому середовищі. Він має вбудовану бібліотеку машинного навчання (MLlib), яка надає інструменти для побудови моделей прогнозування фроду на розподілених обчислювальних кластерах.

Крім того, існують інші середовища програмування, які також можуть бути використані для прогнозування фроду, такі як Julia, Java, C++, і так далі.

Дана кваліфікаційна робота була виконана виключно у середовищі Python.

Основні бібліотеки, які були використані:

- Pandas: ця бібліотека дозволяє працювати з структурованими даними, такими як табличні дані, і забезпечує різноманітні функції для маніпуляції, фільтрації, об'єднання та агрегації даних. Pandas надає багато функціональностей для ефективної роботи з даними, таких як читання і запис даних з різних форматів (наприклад, CSV, Excel, SQL), обробка пропущених значень, обчислення статистичних показників, візуалізація даних та багато іншого.

- NumPy (Numerical Python) - бібліотека, яка надає підтримку для обчислення матриць, векторів та інших багатовимірних масивів. Вона є однією з основних бібліотек для наукових обчислень у Python і надає широкий набір функцій і операторів для роботи з числовими даними. NumPy надає ефективні структури даних для зберігання і маніпуляції числової інформації, включаючи масиви, вектори і матриці. Вона також містить функції для виконання математичних операцій, лінійної алгебри, статистики, обробки зображень та багато іншого.

- Scikit-learn (також відома як sklearn) – ця бібліотека надає широкий набір інструментів і функцій для завдань класифікації, регресії, кластеризації, зменшення розмірності, підбору параметрів та інших завдань, пов'язаних з машинним навчанням.

Scikit-learn побудована на основі інших популярних бібліотек, таких як NumPy та SciPy, і надає простий у використанні та консистентний інтерфейс для роботи з моделями машинного навчання. Вона містить реалізації різних алгоритмів, включаючи методи на основі дерев рішень, метод опорних векторів, навчання з учителем та навчання без учителя, а також інструменти для попередньої обробки даних, оцінки моделей, підбору гіперпараметрів і валідації.

- **Matplotlib та Seaborn:** Matplotlib є основною бібліотекою для створення графіків і візуалізації даних у Python. Вона надає широкий набір функцій для створення різних типів графіків, включаючи лінійні графіки, стовпчасті діаграми, кругові діаграми, розсіювальні графіки, графіки контуру та інші. Matplotlib дозволяє налаштовувати різні аспекти графіків, включаючи заголовки, мітки осей, легенду, колірну палітру та інші елементи, що допомагають передати інформацію візуально.

Seaborn-вищорівнева бібліотека для створення стильних та привабливих графіків у Python. Вона побудована на основі Matplotlib і надає більш простий та зручний інтерфейс для створення графіків зі стандартними налаштуваннями. Seaborn надає спеціалізовані функції для створення графіків, які зазвичай використовуються для візуалізації статистичних даних, таких як графіки розподілу, ящикові діаграми, графіки кореляції та інші. Вона також пропонує красиві стилі оформлення графіків, що допомагають створювати привабливі та професійні візуалізації.

- **CatBoost**—бібліотека для градієнтного бустингу, яка добре працює з категоріальними даними та забезпечує високу точність і стабільність моделей. Назва CatBoost розшифровується як "Categorical Boosting", і її особливість полягає у вбудованій підтримці роботи з категоріальними змінними, що дозволяє зменшити потребу в попередньому перетворенні даних і підвищує якість прогнозування.

- **LightGBM**—високоєфективна бібліотека для градієнтного бустингу, розроблена Microsoft, яка вирізняється швидкістю навчання та низькими вимогами до пам'яті. Бібліотека оптимізована для обробки великих обсягів даних

та має низку особливостей, що забезпечують високу точність і продуктивність, особливо для задач, де важлива швидкість моделі, наприклад, в реальному часі.

Обидві бібліотеки, Matplotlib і Seaborn, дозволяють створювати графіки з великою гнучкістю і налаштуваннями, що дозволяє відтворювати дані у вигляді зрозумілого та інформативного візуального представлення.

Використовуючи мову Python та вищерозглянуті бібліотеки, почнемо виконувати EDA.

2.2 Exploratory Data Analysis

Основним набором даних для дослідження та прогнозування буде датасет «transactions_train.csv», отриманий з сайту Kaggle. У ньому знаходяться дані щодо банківських транзакцій та їх легітимності.

Нижче будуть надані ознаки цього датасету:

Таблиця 2.1

Перелік ознак датасету

Step	відображає одиницю часу в реальному світі. У цьому випадку 1 крок-1 година часу
Type	CASH-IN, CASH-OUT, DEBIT, PAYMENT та TRANSFER.
amount	сума операції в національній валюті.
nameOrig	клієнт, який почав операцію.
oldbalanceOrig	початковий баланс до операції.
newbalanceOrig	баланс клієнта після транзакції.
nameDest	ідентифікатор одержувача транзакції.
oldbalanceDest	початковий баланс одержувача до транзакції.
newbalanceDest	баланс одержувача після транзакції.
isFraud	чи є шахрайством або ні.

Ознайомившись з ознаками даних цього датасету, потрібно завантажити його. Подальші завантаження та відображення ознак та даних датасету будуть виконані завдяки бібліотеці «Pandas».

Виведемо перші три строки датасету для перевірки, чи правильний датасет було завантажено та чи усі ознаки присутні у ньому:

	step	type	amount	nameOrig	oldbalanceOrig	newbalanceOrig	nameDest	oldbalanceDest	newbalanceDest	isFraud
0	1	PAYMENT	9839.64	C1231008815	170136.0	160296.36	M1979787155	0.0	0.0	0
1	1	PAYMENT	1864.28	C1686544295	21249.0	19384.72	M2044282225	0.0	0.0	0
2	1	TRANSFER	181.00	C1305488145	181.0	0.00	C553284065	0.0	0.0	1

Рисунок 2.1-Дані датасету «transactions_train.csv»

Як бачимо, усі ознаки та дані були завантажені коректно. Далі нам потрібно вивести типи ознак, щоб розуміти які ознаки є чисельними та категоріальними:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 6351193 entries, 0 to 6351192
Data columns (total 10 columns):
#   Column          Dtype
---  -
0   step            int64
1   type            object
2   amount          float64
3   nameOrig        object
4   oldbalanceOrig  float64
5   newbalanceOrig  float64
6   nameDest        object
7   oldbalanceDest  float64
8   newbalanceDest  float64
9   isFraud         int64
dtypes: float64(5), int64(2), object(3)
memory usage: 484.6+ MB
```

Рисунок 2.2-Типи ознак

Бачимо, що більшість ознак є числовими, але 3 з них – type, nameOrig та nameDest є категоріальними, тому їх потім потрібно буде за допомогою LabelEncoder закодувати у числові.

Далі розрахуємо кількість пропущених значень (NaN) в кожному стовпці (колонці) датафрейму і повертає ці значення у вигляді масиву (array).

Для кожного стовпця масив містить кількість пропущених значень.

Наприклад, якщо значення для першої колонки (стовпця) рівне 100, то це означає, що в цьому стовпці є 100 пропущених значень.

```
array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0], dtype=int64)
```

Рисунок 2.3-Перевірка пропущених значень

Як бачимо, пропущених значень у цьому датасеті ми не маємо.

Далі перевіримо кількість дубльованих записів у наборі даних train.

```
Out[73]: 0
```

Рисунок 2.4-Кількість дубльованих значень

Дублікатів теж немає, що спрощує аналіз даних.

Тепер нам потрібно вичислити skewness(коефіцієнт асиметрії). В машинному навчанні skewness використовується для оцінки форми розподілу даних та визначення його симетрії або асиметрії. Якщо розподіл даних має високий рівень skewness, це може вказувати на те, що дані не є рівномірно розподіленими та мають виразну симетрію, що може впливати на точність моделі машинного навчання. Отже, для покращення точності моделі машинного навчання може бути корисним провести аналіз skewness і зробити корекцію даних, щоб знизити або усунути асиметрію розподілу даних. Таблиця коефіцієнта асиметрії містить перші 10 стовпців з найвищими значеннями коефіцієнта асиметрії у спадаючому порядку. Коефіцієнт асиметрії вказує на розподіл даних в стовпці: якщо коефіцієнт асиметрії дорівнює нулю, то розподіл є симетричним; якщо він менше нуля, то розподіл має ліву асиметрію, а якщо більше нуля, то праву. Значення коефіцієнта асиметрії від 0,5 і вище вказують на значну асиметрію даних.

Таблиця 2.2

Коефіцієнт асиметрії

	skew
amount	31.050928
isFraud	28.635901
oldbalanceDest	19.934164
newbalanceDest	19.362310
oldbalanceOrig	5.243790
newbalanceOrig	5.172421
step	0.338249

Як бачимо, коефіцієнт асиметрії дуже великий, його потрібно відкоригувати перед тим, як робити модель прогнозування.

Для цього спочатку, для кращого розуміння різниці між невідкоригованими даними та відкоригованими, відобразимо графіки нинішніх даних:

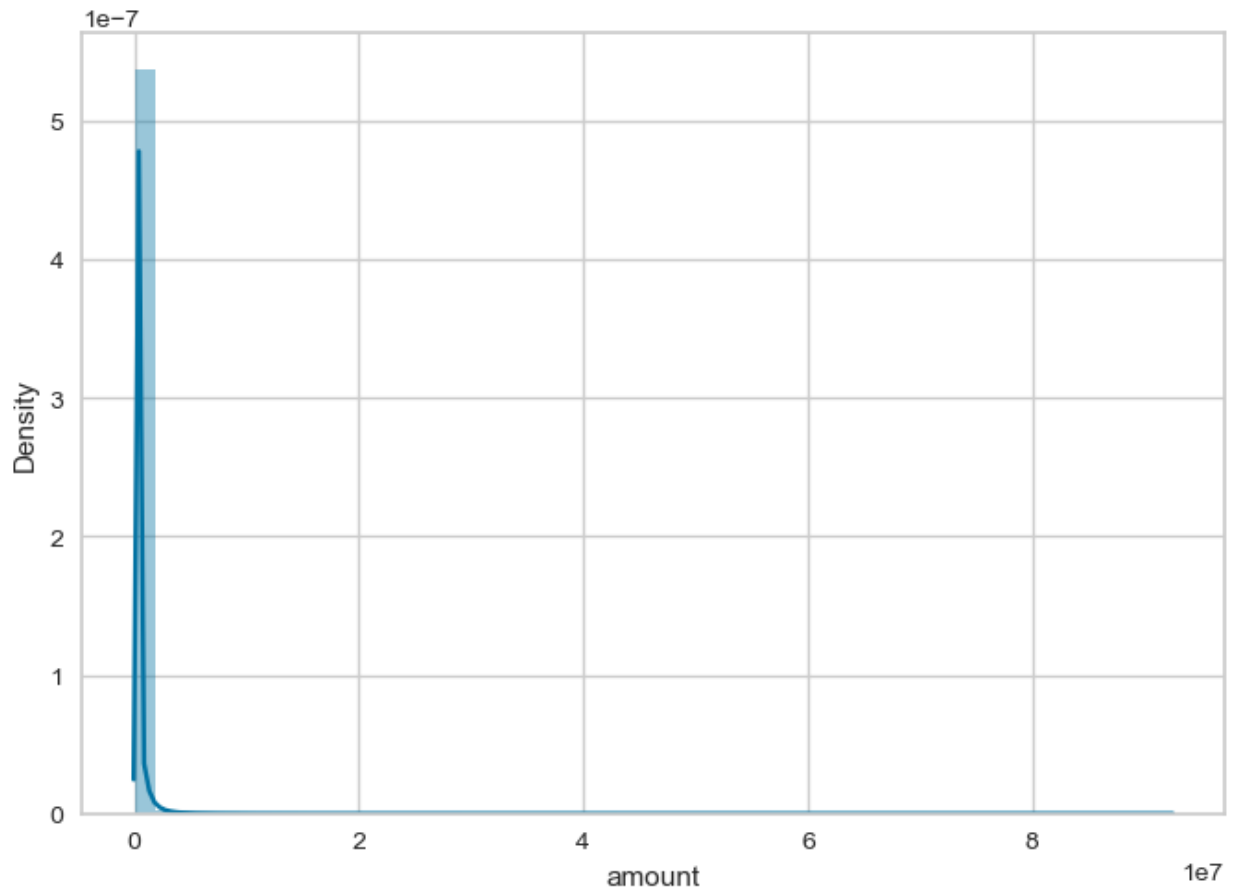


Рисунок 2.5-Відображення асиметрії даних

Виконаємо корекцію розподілу значень стовпця amount за допомогою логарифмування:

КА до логарифмування: 31.050928455018084
КА після логарифмування: -0.5549658313745637

Рисунок 2.6-Відкореговані дані

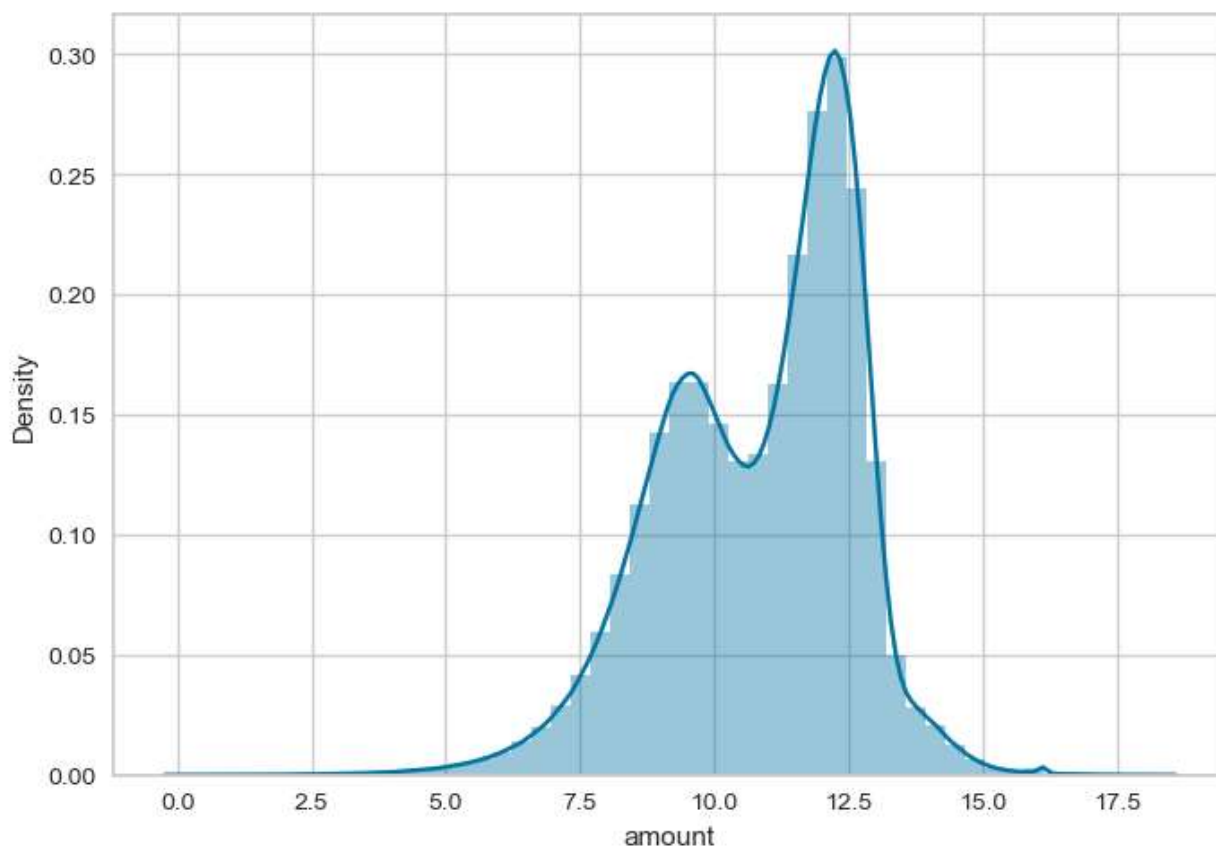


Рисунок 2.7-Графік після логарифмування

Отже, після логарифмування, значення коефіцієнта асиметрії зменшилось, що означає, що розподіл став менш скошеним та більш симетричним.

Далі обчислимо коефіцієнт ексцесу. Обчислення куртозису для кожного стовпця дозволяє оцінити, наскільки вибірка відрізняється від нормального розподілу. Значення куртозису можуть бути позитивними або від'ємними.

Таблиця 2.3

Коефіцієнти ексцесу

	Kurtosis
step	0.246047
amount	1803.410673
oldbalanceOrig	32.875430
newbalanceOrig	32.003795
oldbalanceDest	950.015902
newbalanceDest	863.076045
isFraud	818.015079

Значення куртозису більше 0 вказує на більш "концентрований" розподіл з тяжкими хвостами, тоді як значення менше 0 вказує на менш "концентрований" розподіл з меншими хвостами. Значення навколо 0 вказують на близькість до нормального розподілу.

1. Для стовпця "step" куртозис дорівнює 0.246047, що вказує на наближеність до нормального розподілу.

2. Для стовпця "amount" куртозис дорівнює 1803.410673, що вказує на дуже велику хвостатість розподілу. Це означає, що значення "amount" розподілені широко з великими хвостами, що може вказувати на наявність викидів або аномалій у даних.

3. Для стовпців "oldbalanceOrig", "newbalanceOrig", "oldbalanceDest" і "newbalanceDest" куртозиси дорівнюють відносно невеликим значенням, що вказує на наближеність до нормального розподілу. Це означає, що значення в цих стовпцях мають меншу хвостатість і більшу концентрацію навколо середнього значення.

4. Для стовпця "isFraud" куртозис дорівнює 818.015079, що вказує на велику хвостатість розподілу.

Тепер обчислемо дисперсії для всіх стовпчиків. Аналіз дисперсії допомагає зрозуміти варіативність та розподіл даних у відповідних стовпцях вашого датафрейму. Ця інформація може бути використана для подальшого аналізу, видалення аномалій, вибору ознак для моделювання та виконання інших завдань аналізу даних.

Таблиця 2.4

Дисперсія стовбців

	var
isFraud	1.213571e-03
step	1.990008e+04
amount	3.643704e+11
oldbalanceOrig	8.351864e+12
newbalanceOrig	8.561904e+12
oldbalanceDest	1.155268e+13
newbalanceDest	1.350043e+13

Зі значень дисперсії можна зробити наступні спостереження:

1. isFraud: стовпець, що вказує на наявність або відсутність шахрайства. Значення дисперсії дуже низькі (1.213571e-03), що може свідчити про те, що дані в цьому стовпці майже однакові або майже не змінюються. Відсутність великої варіативності в цьому стовпці може бути пов'язана з тим, що більшість транзакцій не є шахрайськими.

2. step: стовпець, який вказує на часовий крок транзакції. Значення дисперсії досить великі (1.990008e+04), що свідчить про значну різноманітність значень і можливо про довгий часовий проміжок між окремими транзакціями.

3. amount, oldbalanceOrig, newbalanceOrig, oldbalanceDest, newbalanceDest: Ці стовпці пов'язані з сумами грошей або балансами на рахунках. Значення дисперсії у всіх цих стовпцях дуже великі, що свідчить про значну варіативність у сумах та балансах. Це може бути пов'язано з різними типами транзакцій та рахунків, які мають різні розміри та зміни у значеннях.

Зі спостережень дисперсій можемо зробити висновок, що аномалій не знайдено.

Далі обчислемо IQR або міжквартильний діапазон. Обчислення IQR дозволяє визначити "нормальний" діапазон значень у стовпці. Зазвичай вважається, що значення, що виходять за цей діапазон, можуть бути потенційними викидами або аномаліями.

Таблиця 2.5

IQR

step	179.00
amount	195326.90
oldbalanceOrig	107346.00
newbalanceOrig	144365.15
oldbalanceDest	943866.12
newbalanceDest	1112791.08
isFraud	0.00
dtype: float64	

На підставі отриманих значень IQR для кожного стовпця можна зробити наступні висновки:

1. step: Розмах значень стовпця step становить 179, що означає, що дані мають відносно невелику варіабельність в цьому стовпці.

2. amount: Розмах значень стовпця amount досить великий і становить 195,326.90. Це свідчить про значну варіативність суми транзакцій.

3. oldbalanceOrig та newbalanceOrig: Ці стовпці також мають значну варіабельність, оскільки їх розмахи становлять відповідно 107,346.00 і 144,365.15. Це означає, що вихідні та нові баланси походять з різних діапазонів значень.

4. oldbalanceDest та newbalanceDest: Ці стовпці також мають значну варіабельність, оскільки їх розмахи становлять відповідно 943,866.12 і 1,112,791.08. Це означає, що старі та нові баланси отримувачів також знаходяться в широкому діапазоні значень.

Ми можемо видалити викиди із датасету за допомогою міжквантильного діапазону, але це може видалити значення з шахрайських транзакцій, що може зменшити шанси на виявлення фроду під час прогнозування. Але ми можемо

використовувати алгоритми, стійкі до викидів, такі як дерева рішень, випадковий ліс або алгоритм XGBoost, що допоможе нам не втратити значень для прогнозування фроду.

Далі створимо pair plot, який використовує бібліотеку Seaborn для візуалізації взаємозв'язків між змінними з датасету з перших 50 000 рядків.

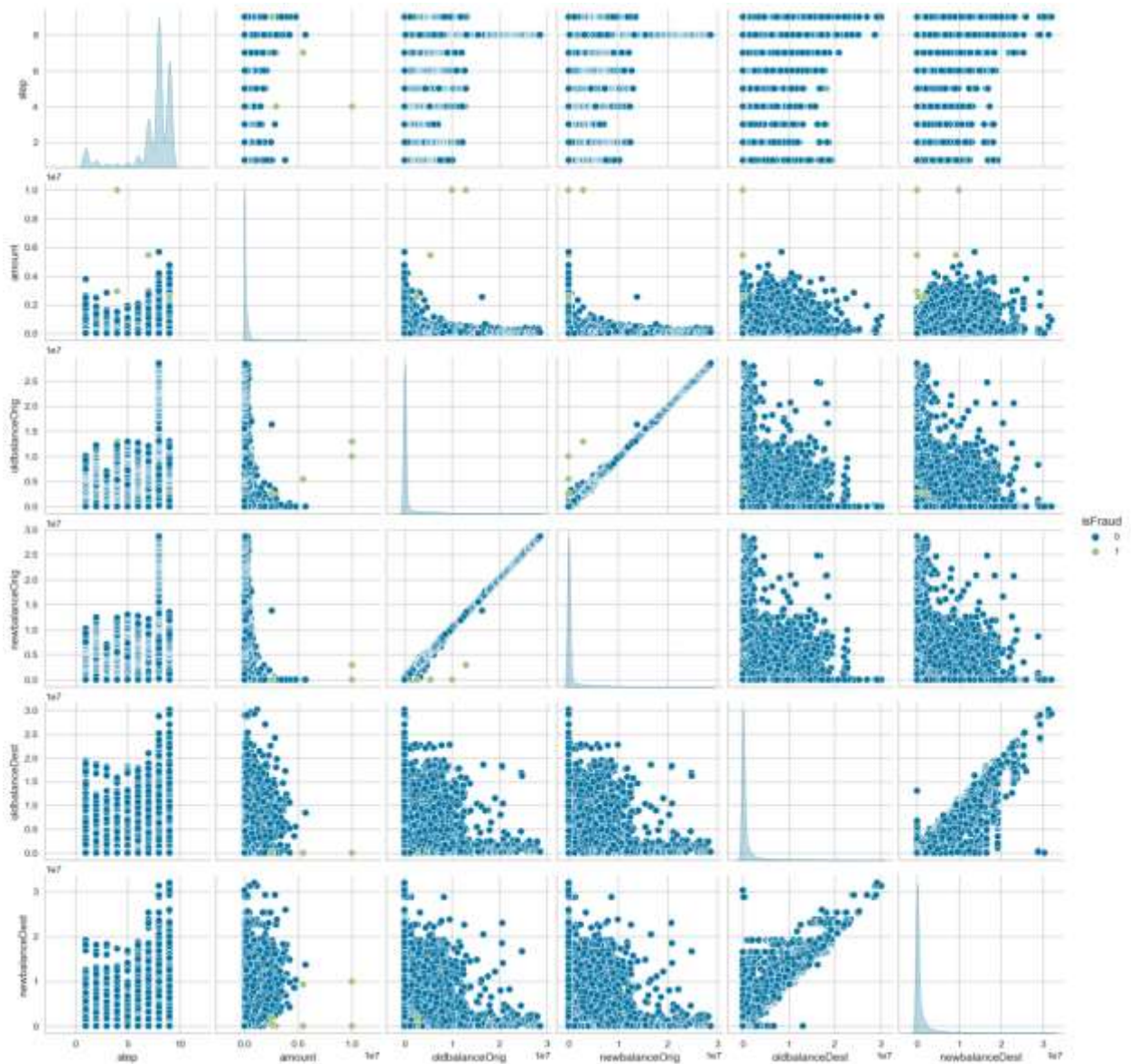


Рисунок 2.8-Графік взаємозв'язків ознак датасету до логарифмування

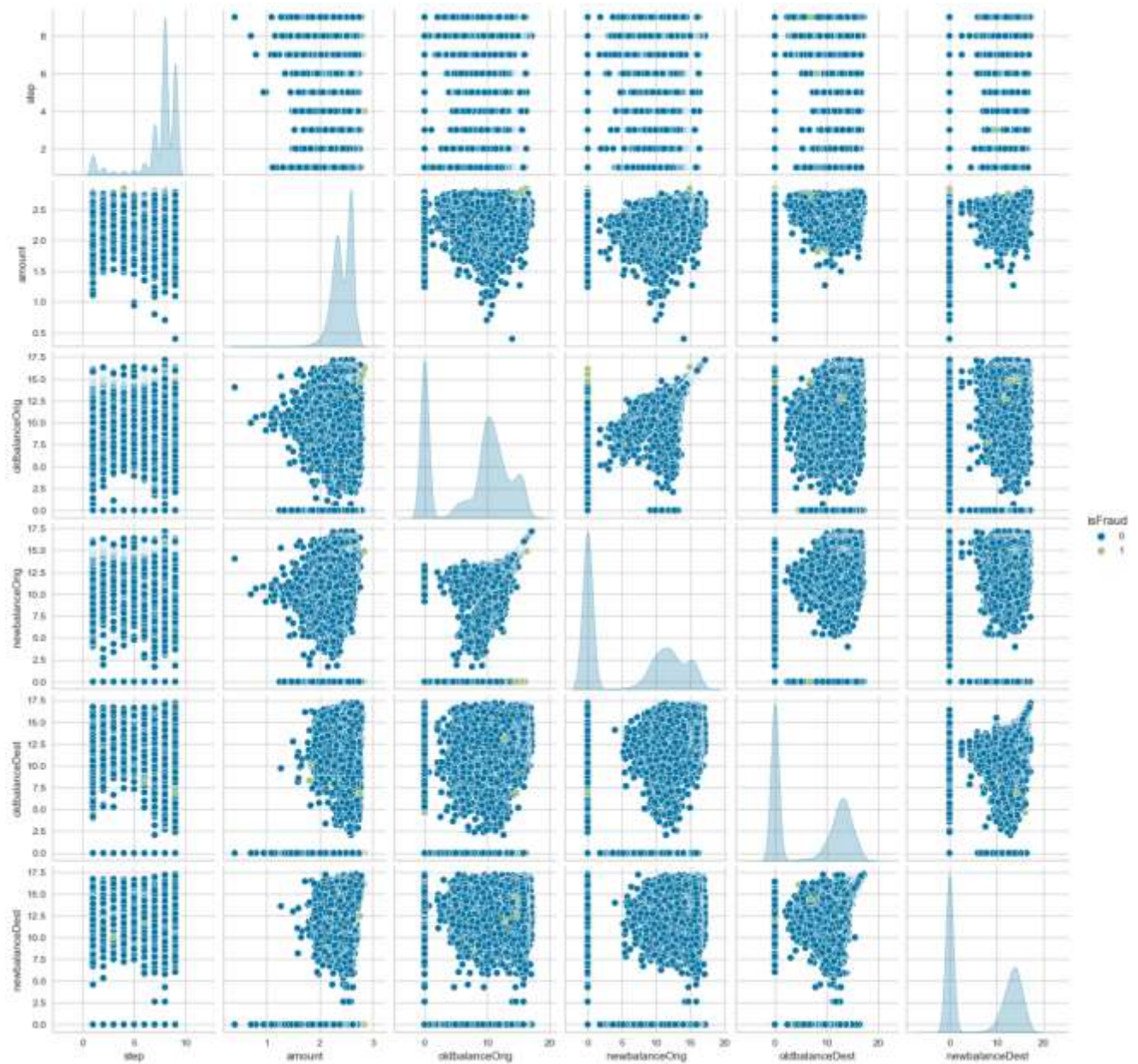


Рисунок 2.9-Графік взаємозв'язків ознак датасету після логарифмування

Як можемо побачити, `newbalanceOrig` та `oldbalanceOrig` мали сильні взаємозв'язки до логарифмування, також значимий зв'язок мали `newbalanceDest` та `oldbalanceDest`. Сильний взаємозв'язок між цими змінними може бути важливим для аналізу фінансових транзакцій, наприклад, це може свідчити про здійснення переказів коштів між різними рахунками або зміну балансу після здійснення операцій, але через логарифмування ми втратили деяку частину зв'язку, тому кращим варіантом буде логарифмувати лише ознаку `amount`, бо вона має найбільший коефіцієнт асиметрії.

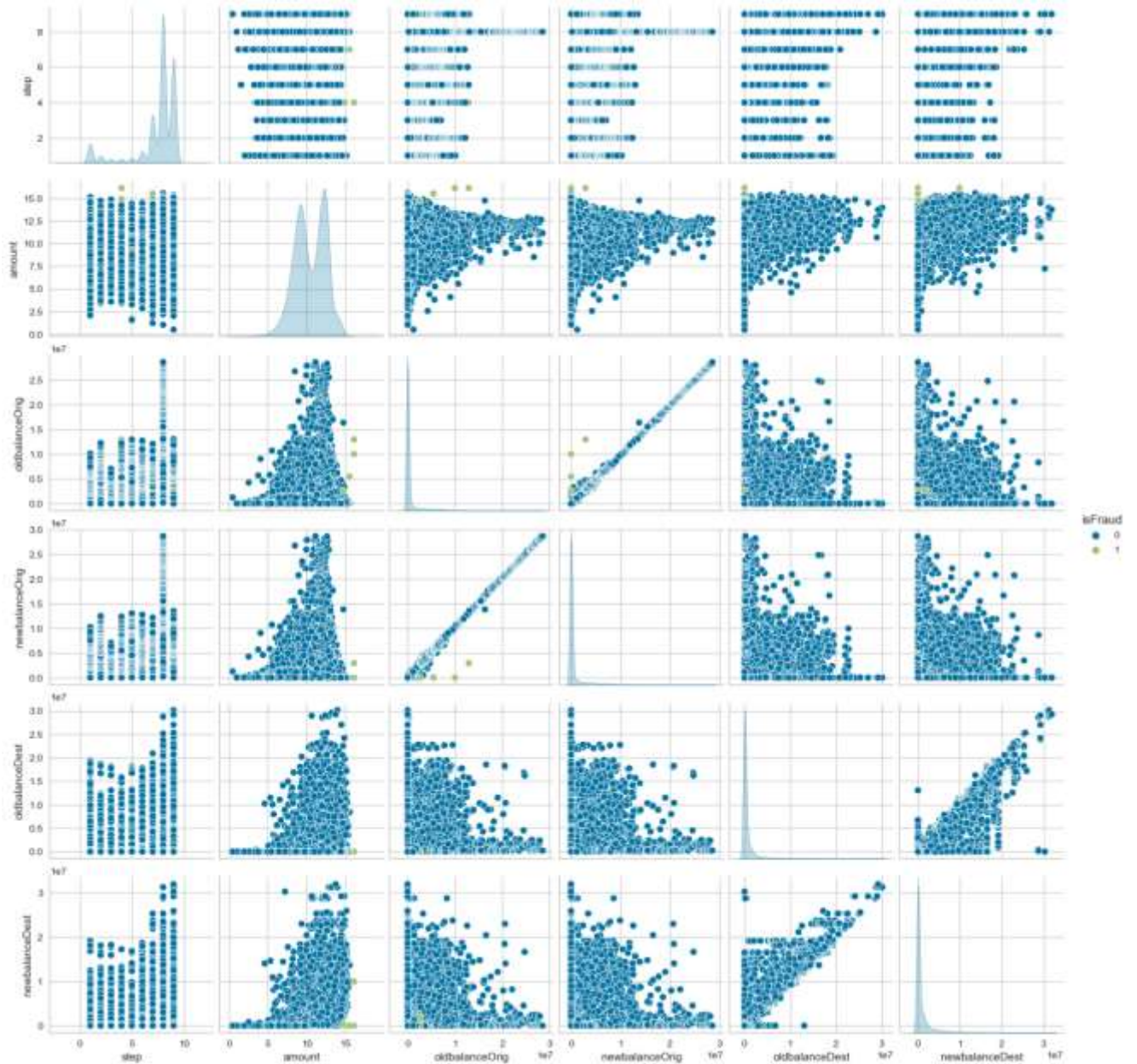


Рисунок 2.10-Графік взаємозв'язків ознак датасету після логарифмування лише ознаки amount

Далі створимо кореляційну матрицю для датасету і візуалізує її в графічному вигляді за допомогою функції `corr_plot` з бібліотеки `klib`.

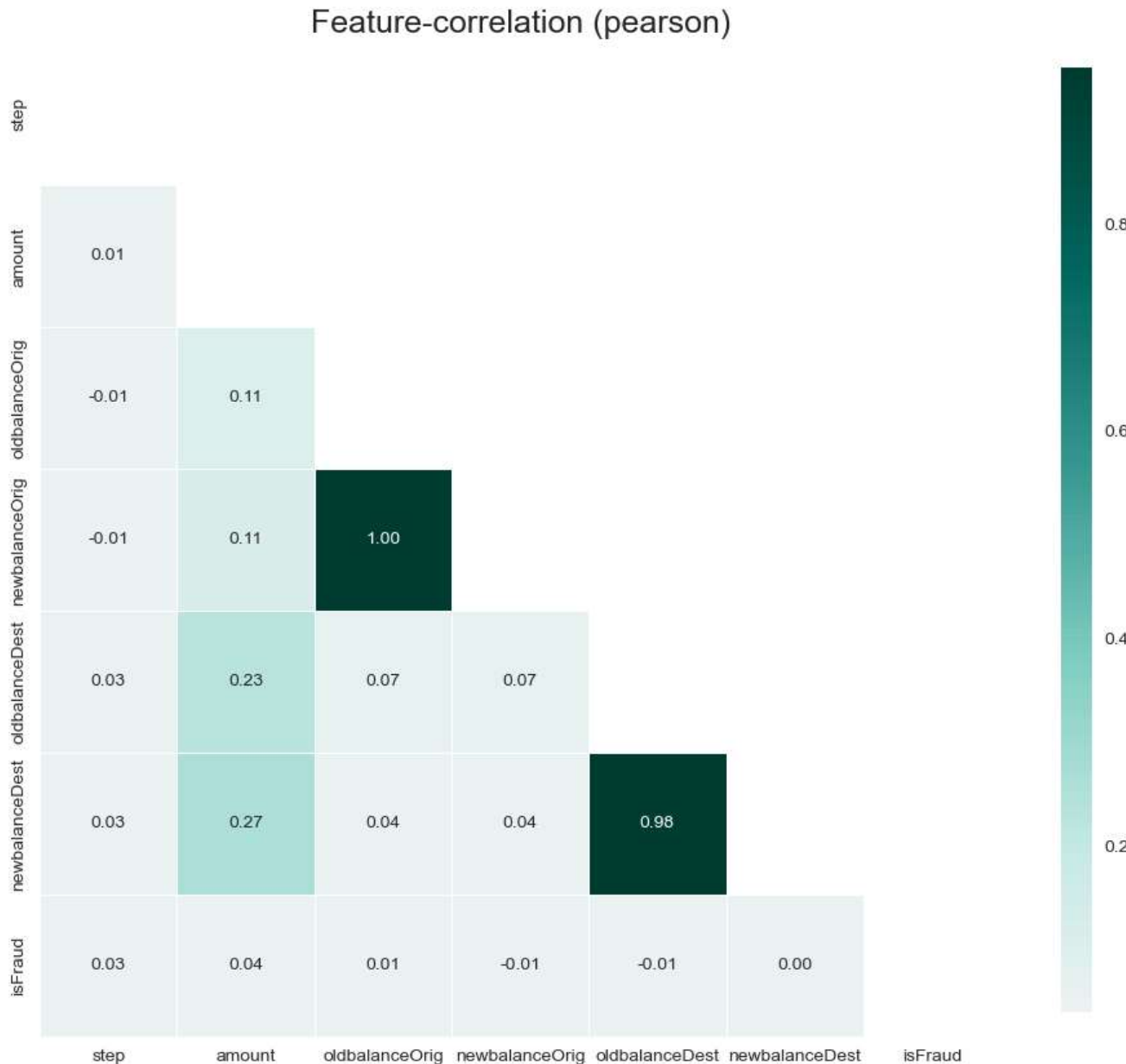


Рисунок 2.11-Графік кореляції Пірсона

Дійсно бачимо, взаємозв'язок `newbalanceOrig` та `oldbalanceOrig` є майже 100%, що означає що кожна зміна даних у однієї з цих двох ознак прямо впливає на іншу.

Далі змінимо тип даних колонок у датафсеті на числові та категорійні за допомогою методу `.apply()` з бібліотеки `Pandas`. Це може бути корисно, якщо деякі стовпці були помилково розпізнані як рядкові об'єкти або категорії, і їх потрібно перетворити на числові типи, щоб проводити числові операції і аналіз.

Далі створимо змінні x та y , що містять вхідні дані та відповідні мітки класів для задачі класифікації.

Далі створимо дві змінні: `cat_columns` та `num_columns`. Змінна `cat_columns` містить назви усіх категорійних колонок у X , а змінна `num_columns` містить назви усіх числових колонок. Ці змінні можуть бути використані для подальшої обробки та аналізу даних відповідного типу.

Далі створимо 6 гістограм, по одній для кожної змінної з числовими значеннями (що містяться в змінній `num_columns`) з даних навчального набору даних. Кожна гістограма містить розподіл значень змінної, показаний за допомогою кількох стовпчиків (бінів), де кожен стовпчик відповідає певному діапазону значень змінної. Також на гістограмі відображений неперервний розподіл щільності (`kde`), який дозволяє оцінити швидкість зміни частоти виникнення значень. У кожній гістограмі мітки на осі x показують значення змінної, а на осі y - частоту виникнення значень.

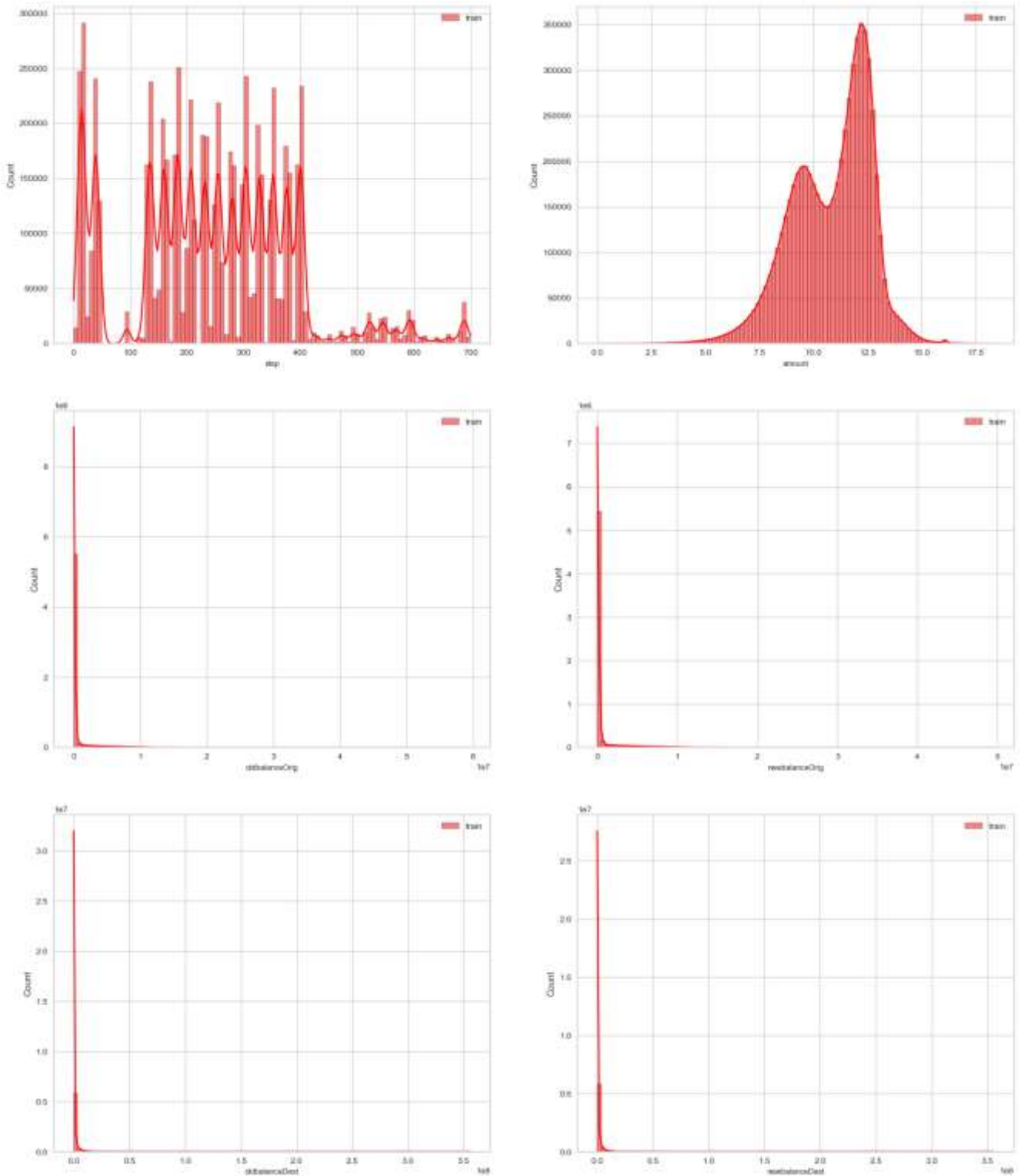


Рисунок 2.12-Гістограми розподілів змінної

Далі ще раз обчислимо кореляцію між числовими ознаками у наборі даних, застосовуючи градієнтний колір для більш наочного показу кореляційних коефіцієнтів :

	step	amount	oldbalanceOrig	newbalanceOrig	oldbalanceDest	newbalanceDest	isFraud
step	1.000000	0.007333	-0.009113	-0.009201	0.028303	0.026508	0.025495
amount	0.007333	1.000000	0.106908	0.111480	0.228012	0.266220	0.039389
oldbalanceOrig	-0.009113	0.106908	1.000000	0.998857	0.066301	0.042019	0.009226
newbalanceOrig	-0.009201	0.111480	0.998857	1.000000	0.067852	0.041853	-0.008322
oldbalanceDest	0.028303	0.228012	0.066301	0.067852	1.000000	0.976550	-0.005657
newbalanceDest	0.026508	0.266220	0.042019	0.041853	0.976550	1.000000	0.000496
isFraud	0.025495	0.039389	0.009226	-0.008322	-0.005657	0.000496	1.000000

Рисунок 2.13-Кореляція між відкоригованими даними ознак

Можемо побачити, що тепер newbalanceDest та oldbalanceDest мають сильний зв'язок між собою.

Побудуємо стовпчасту діаграму, яка показує кореляцію між змінною isFraud (цільова змінна) та всіма іншими числовими змінними в наборі даних.

Зверху вниз, змінні розташовані у порядку спадання їх кореляції зі змінною isFraud:

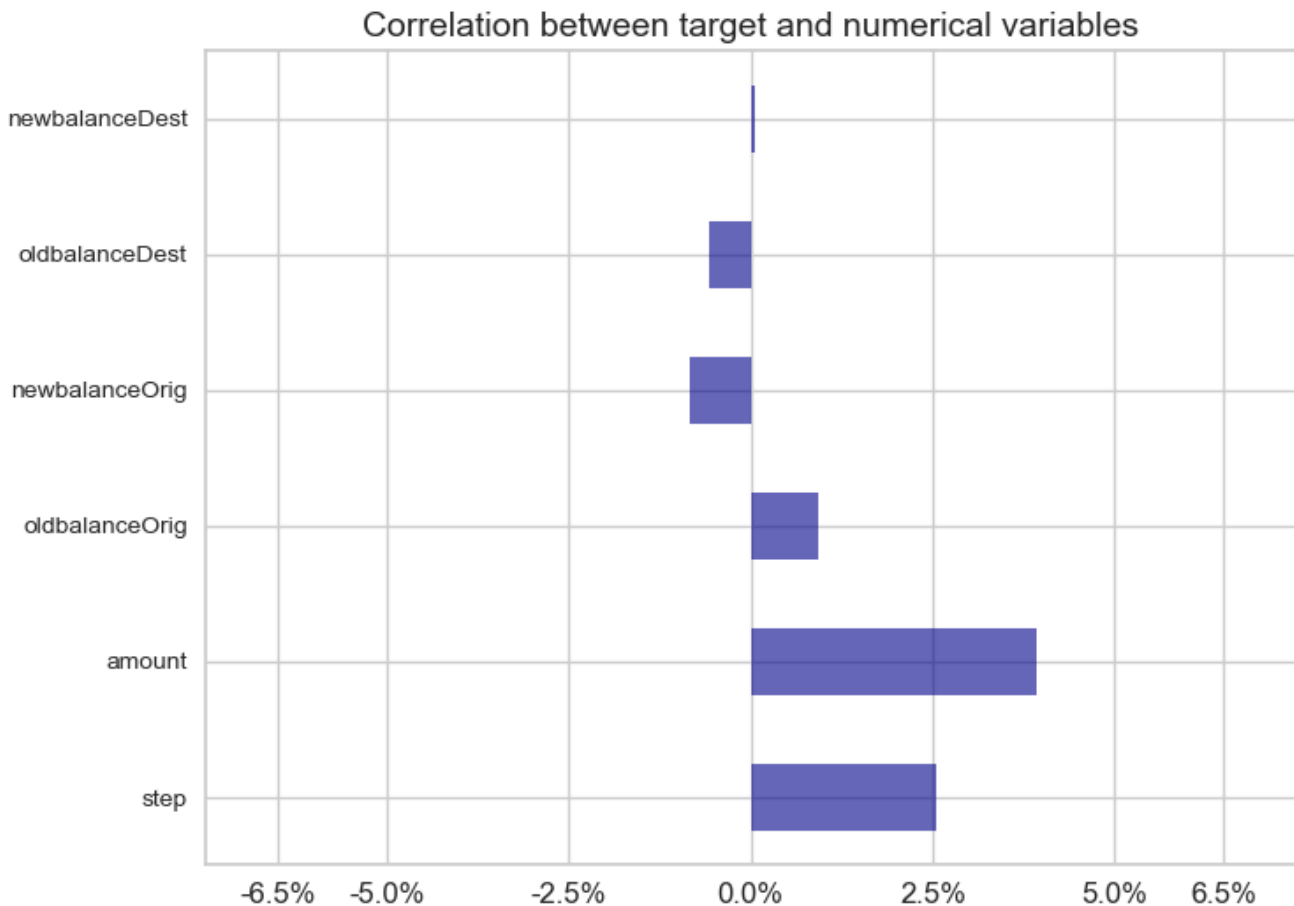


Рисунок 2.14-Кореляція між isFraud та числовими змінними

Від'ємне значення кореляції означає, що змінні рухаються в протилежних напрямках, тоді як додатне значення означає, що змінні рухаються в одному напрямку.

Графік чітко відображає, що amount найбільш зв'язаний з isFraud.

Побудуємо стовпчикову діаграму, яка показує кількість унікальних значень для кожної категоріальної ознаки з набору даних, де:

у - кількість унікальних значень кожної категоріальної ознаки з набору даних `train[cat_columns]`.

х - імена категоріальних ознак з `train[cat_columns]`.

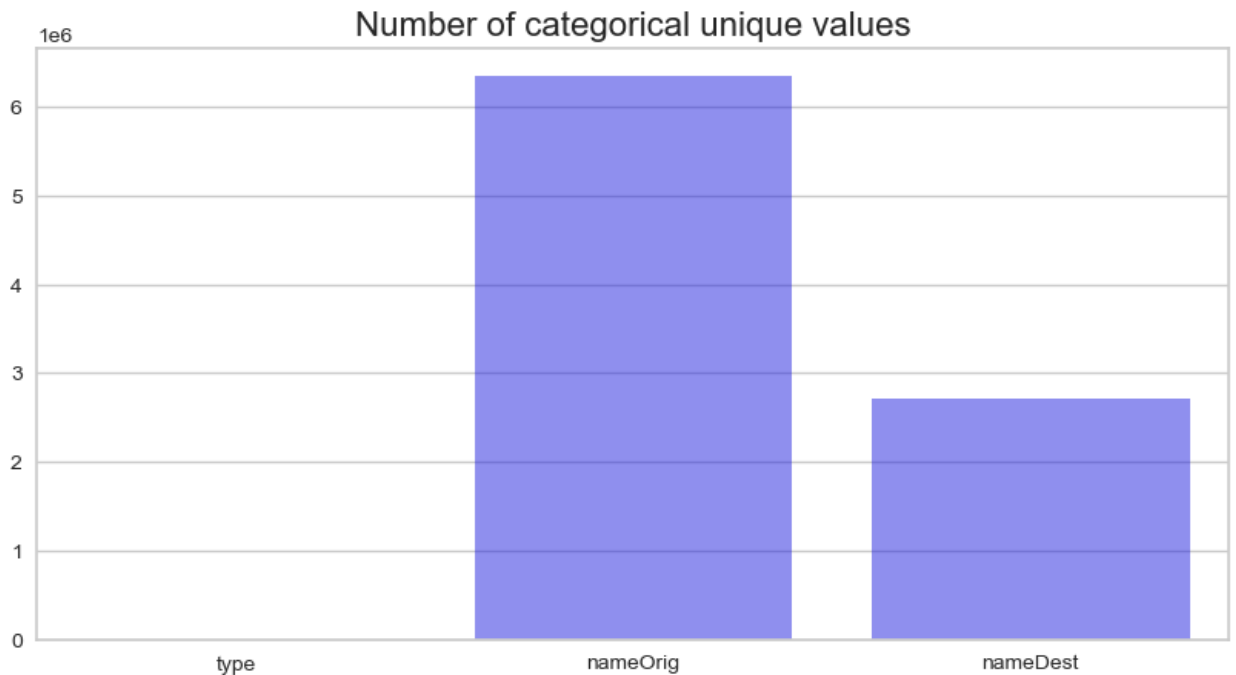


Рисунок 2.15-Кількість унікальних значень у датасеті

Інформація про унікальні значення в категоріальних стовпцях набору даних може бути корисною з декількох причин:

1. Розуміння розподілу категорій: Знання кількості унікальних значень допомагає нам отримати уявлення про розподіл категорій у стовпці. Це може бути важливою інформацією для подальшого аналізу та врахування особливостей розподілу у моделюванні.

2. Виявлення високовимірних стовпців: Якщо категоріальний стовпець має велику кількість унікальних значень, це може свідчити про високу вимірність цього стовпця. Високовимірні стовпці можуть бути складними для обробки та аналізу, а також можуть призводити до перенавчання моделей. Знання про кількість унікальних значень допоможе нам виявити такі стовпці та прийняти рішення щодо їх обробки або включення в модель.

3. Кодування категоріальних змінних: при роботі з категоріальними змінними для моделювання часто потрібно закодувати їх у числовий формат. Знання про кількість унікальних значень допомагає визначити, який метод кодування буде найкращим в конкретному випадку. Наприклад, якщо категоріальний стовпець має декілька унікальних значень, можна розглянути

застосування методу "one-hot encoding", тоді як для стовпця з великою кількістю унікальних значень може бути доцільним використовувати кодування за допомогою "label encoding" або "target encoding".

Як бачимо, велика кількість унікальних категоріальних значень у цьому датасеті потребує кодування для того, щоб наша модель класифікації працювала без помилок.

Обчислємо кількість записів з кожним значенням ознаки "isFraud" у наборі даних.:

```
0      6343476
1       7717
Name: isFraud, dtype: int64
```

Рисунок 2.16-Кількість шахрайських та легітимних транзакцій

Це означає, що в наборі даних є 6343476 транзакцій, які не є шахрайськими (значення 0) та 7717 транзакцій, які є шахрайськими (значення 1) за ознакою "isFraud". Тому можемо підтвердити, що наш набір даних є незбалансованим.

Відобразимо на графіку цю незбалансованість для розуміння степені незбалансованості:

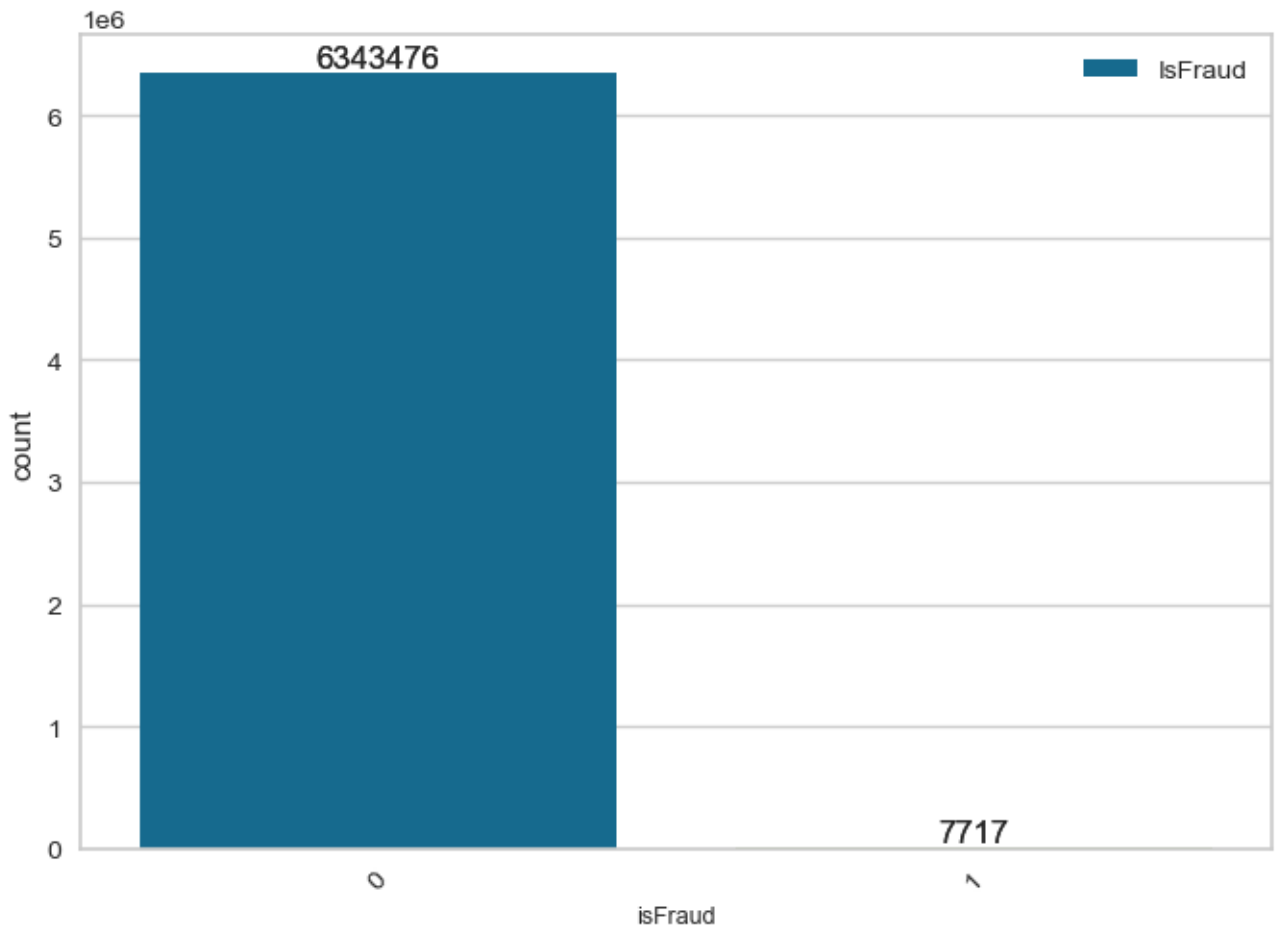


Рисунок 2.17-Графік незбалансованості даних шахрайств

Побудуємо графіки скрипичних діаграм (violinplot) для візуалізації зв'язку між числовими ознаками та цільовою змінною 'isFraud'.

Побудова цих графіків має такі цілі:

1. Візуалізація розподілу значень числових ознак залежно від класу 'IsFraud' (шахрайство або не шахрайство). Це допомагає зрозуміти, які значення ознак мають різниці або схожості між класами. Графіки скрипки дозволяють побачити щільність значень ознак, медіану, міжквартильний розмах та діапазон значень для кожного класу.

2. Виявлення потенційних розбіжностей у розподілі ознак між класами 'IsFraud'. Якщо графіки скрипки для певної ознаки мають різні форми або розташування медіани, це може вказувати на важливість цієї ознаки для класифікації шахрайства.

3. Виявлення викидів або незвичайних значень у ознаках для кожного класу. Графіки скрипки можуть допомогти виявити аномалії або екстремальні значення у розподілі ознак, які можуть бути пов'язані з шахрайською активністю.

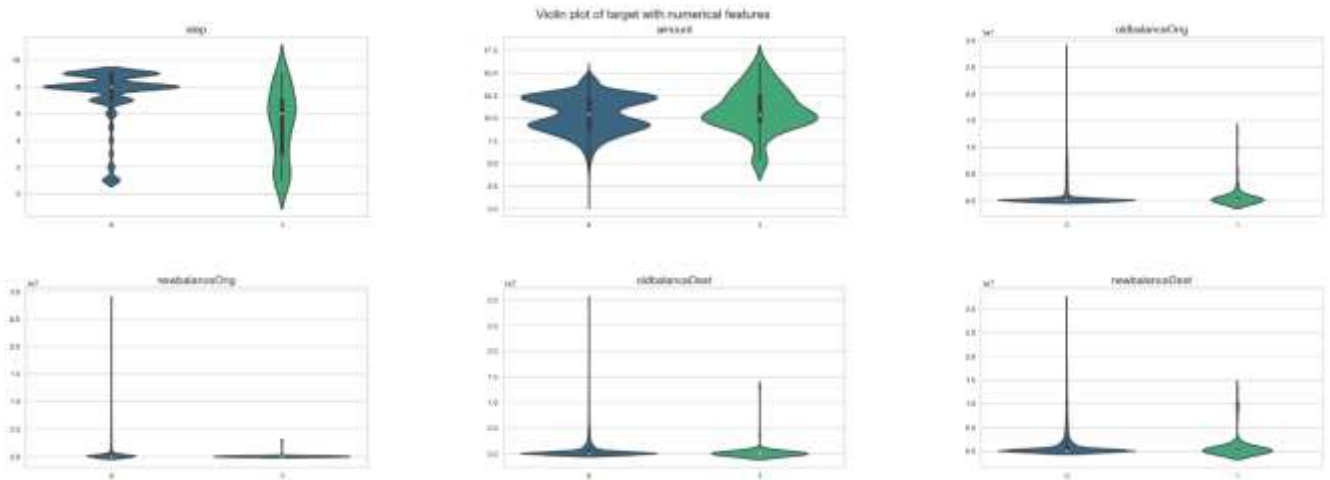


Рисунок 2.18-Графіки скрипичних діаграм

Побудова цих графіків має такі цілі:

1. Візуалізація розподілу значень числових ознак залежно від класу 'IsFraud' (шахрайство або не шахрайство). Це допомагає зрозуміти, які значення ознак мають різниці або схожості між класами. Графіки скрипки дозволяють побачити щільність значень ознак, медіану, міжквартильний розмах та діапазон значень для кожного класу.

2. Виявлення потенційних розбіжностей у розподілі ознак між класами 'IsFraud'. Якщо графіки скрипки для певної ознаки мають різні форми або розташування медіани, це може вказувати на важливість цієї ознаки для класифікації шахрайства.

3. Виявлення викидів або незвичайних значень у ознаках для кожного класу. Графіки скрипки можуть допомогти виявити аномалії або екстремальні значення у розподілі ознак, які можуть бути пов'язані з шахрайською активністю.

Можемо побачити, що найбільший розмах ширини значень `isFraud` має `newbalanceOrig`, також бачимо підтвердження що ознаки `newbalanceDest` та `oldbalanceDest` мають сильний зв'язок між собою, бо їх розмах майже ідентичний.

Створимо графік із 6 підграфіками, кожен з яких показує KDE (ядерну оцінку щільності) числової ознаки у датафреймі.

Ядерна оцінка щільності є корисним інструментом для візуалізації та аналізу розподілу даних, виявлення відмінностей між різними групами та оцінки ймовірності.

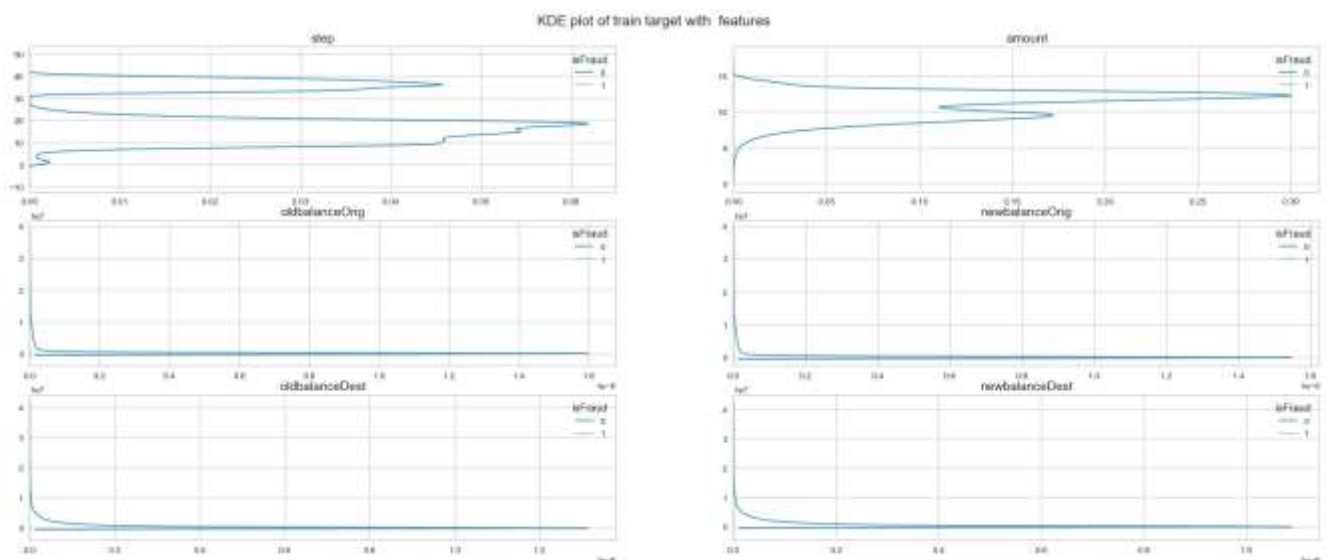


Рисунок 2.19-Ядерна оцінка

Трансформуємо дані з використанням алгоритму TSNE (t-distributed stochastic neighbor embedding) для зменшення кількості вимірів даних і візуалізації їх в двовимірному просторі. Спочатку з набору даних вилучаються нечислові ознаки, які не можуть бути використані для tsne-преобразування. Після цього з застосуванням `RobustScaler` числові дані нормалізуються і підготовлюються для використання в tsne. Нарешті, виконується само tsne-преобразування, яке зменшує кількість вимірів з метою візуалізації даних в двовимірному просторі. Результатом є numpy масив з двома стовпцями, що представляють нові координати кожної точки даних в двовимірному просторі.

Потім виконаємо візуалізація точок на площині x-y залежно від значення стовпця isFraud:

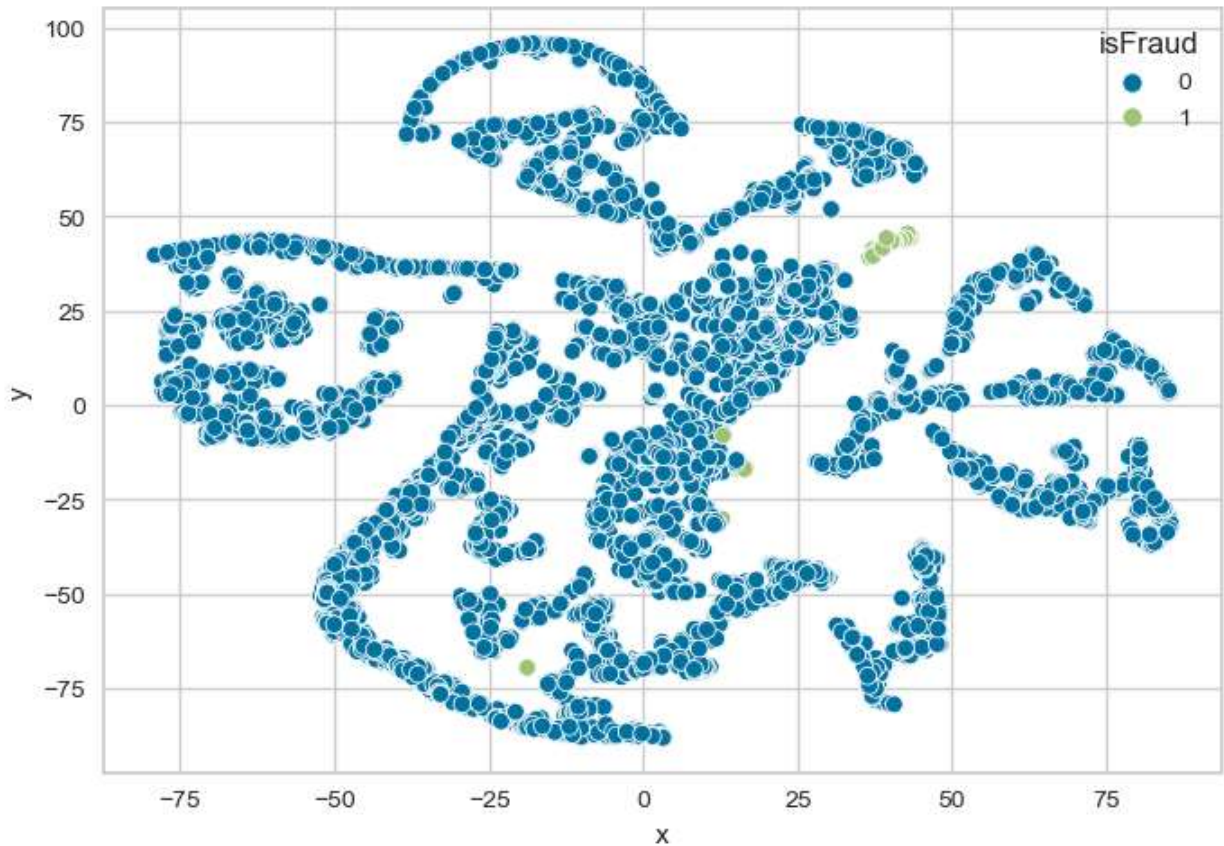


Рисунок 2.20-Візуалізація точок шахрайства

Як ми бачимо, нелінійне перетворення дає нам деякий кластер шахрайства, добре відокремлений від нешахрайських даних.

Далі розіб'ємо навчальний датасет на дві частини - навчальний та тестовий. Розмір тестового датасету встановлюється на 10% від загального розміру. Крім того, цей метод робить збалансований поділ даних, використовуючи параметр `stratify=y`. Це означає, що кожен окремий клас має представника як у навчальному, так і в тестовому датасетах.

```
635120 rows in test set vs. 5716073 in training set. 9 Features.
```

Рисунок 2.21-Результат розділення

На цьому етапі розвідковий аналіз даних завершено. У наступному розділі приступимо до програмування та тесту моделей прогнозування шахрайства.

2.4 Побудова моделей прогнозування

Перед тим як будувати модель XGBoost, спочатку кластеризуємо навчальний датасет за допомогою двох методів: K-means та Fuzzy Clustering. Ці методи кластеризації можуть поліпшити якість нашої моделі, ізолюючи більшість шахрайств у одному кластері.

2.4.1 Модель прогнозування з кластеризацією K-means

Перед тим як кластеризувати дані, спочатку закодуємо категоріальні ознаки type, nameOrig та nameDest за допомогою Label Encoder.

```

step          int64
type          int32
amount        float64
nameOrig      int32
oldbalanceOrig float64
newbalanceOrig float64
nameDest      int32
oldbalanceDest float64
newbalanceDest float64
dtype: object
step          int64
type          int32
amount        float64
nameOrig      int32
oldbalanceOrig float64
newbalanceOrig float64
nameDest      int32
oldbalanceDest float64
newbalanceDest float64
dtype: object

```

Рисунок 2.22-Результати кодування

Далі нам потрібно нормалізувати дані за допомогою алгоритму Standard Scaler.

Після цього нам потрібно знати, на скільки кластерів потрібно розділити дані. Для цього реалізуємо метод "Лікоть" (Elbow Method) для визначення оптимальної кількості кластерів у алгоритмі K-means. Він обчислює значення інерції (суми квадратів відстаней) для різних кількостей кластерів і показує графік "крива лікоть". Оптимальне число кластерів визначається за допомогою знаходження точки "лікоть" на графіку.

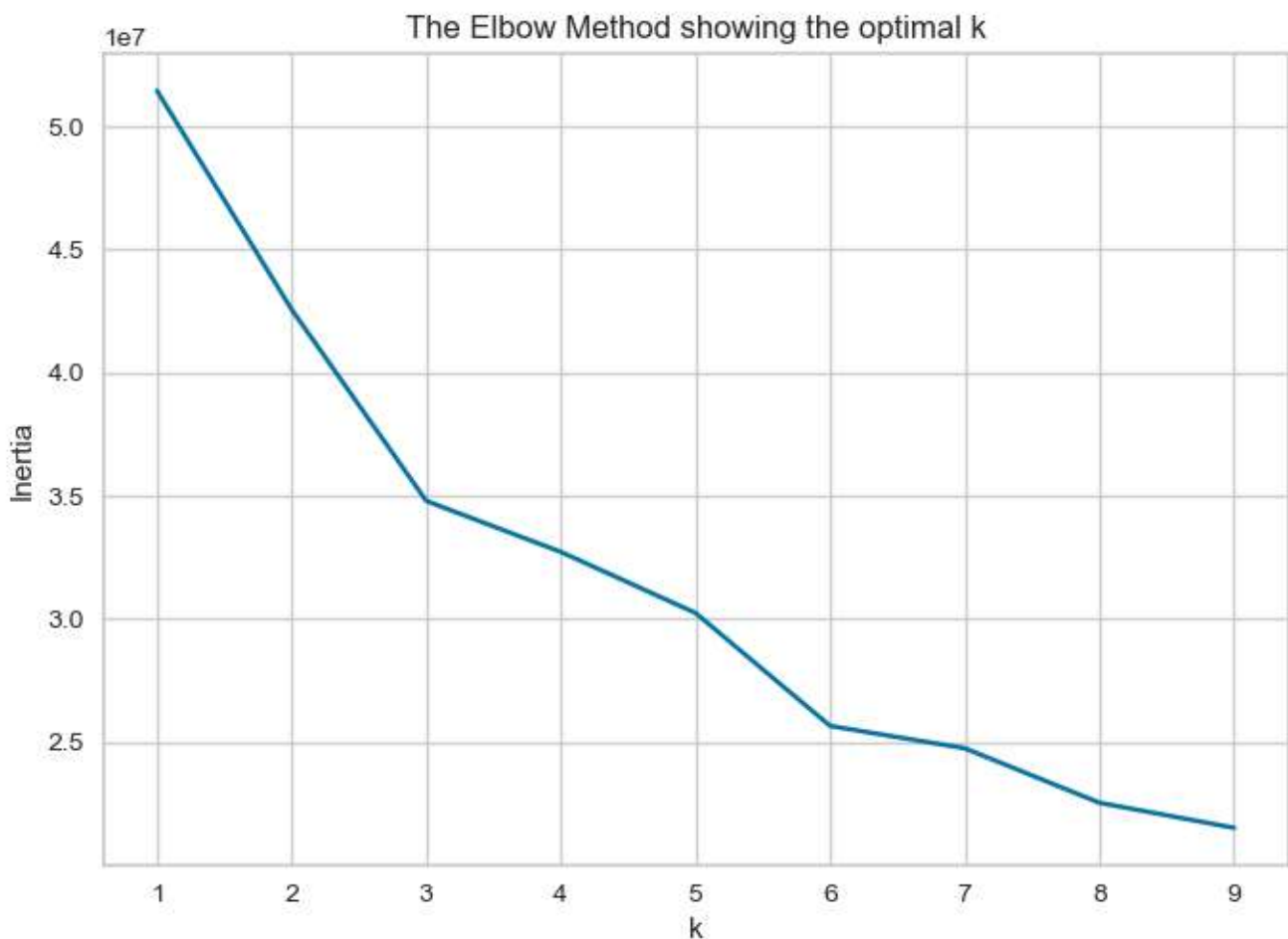


Рисунок 2.23-Оптимальна кількість кластерів

Після цього реалізуємо метод K-means.

Виконуємо наступні кроки:

1. Застосовується кодування міток класів у вихідних даних за допомогою LabelEncoder. Вихідні мітки класів у `y_train` перетворюються на числові значення.
2. Використовується отримане оптимальне число кластерів (елбоу індекс) для створення інстанції алгоритму кластеризації K-means (MiniBatchKMeans).

3. Алгоритм K-means навчається на навчальних даних (`X_train_scaled`) і призначає кластери кожній точці даних у навчальному наборі.

Отримані кластерні мітки для навчального набору зберігаються в змінній `cluster_labels_train`.

Далі перевіримо в якому кластері найбільша кількість шахрайств:

Таблиця 2.6

Процент шахрайств у кластерах

isFraud	0	1	All
cluster			
0	0.545078	0.518790	0.545046
1	0.413195	0.468539	0.413283
2	0.041727	0.012871	0.041692
All	1.000000	1.000000	1.000000

Як бачимо, таблиця показує, що у кластері №0 найбільша кількість шахрайств.

Перевіримо також який процент шахрайств є у кожному з кластерів:

```
Fraud percentage in Cluster 0: 56.21761658031088
Fraud percentage in Cluster 1: 42.35751295336787
Fraud percentage in Cluster 2: 1.4248704663212435
100.0
```

Рисунок 2.24-Процент шахрайств у кожному з кластерів

Бачимо, що що кластер №0 дійсно є лідером у проценті шахрайств. Також ми маємо біль детальну інформацію щодо інших кластерів, за допомогою якої ми можемо зрозуміти, хто з останніх двох кластерів має більшу кількість шахрайств, а саме кластер № 1. У такому випадку, можемо провести додатковий експеримент з підвищення якості моделі прогнозування – скомбінувати кластер №0 та №1 та зробити на цих даних модель прогнозування.

Перевіримо точну кількість для перевірки, що програма правильно обчислила процент шахрайств:

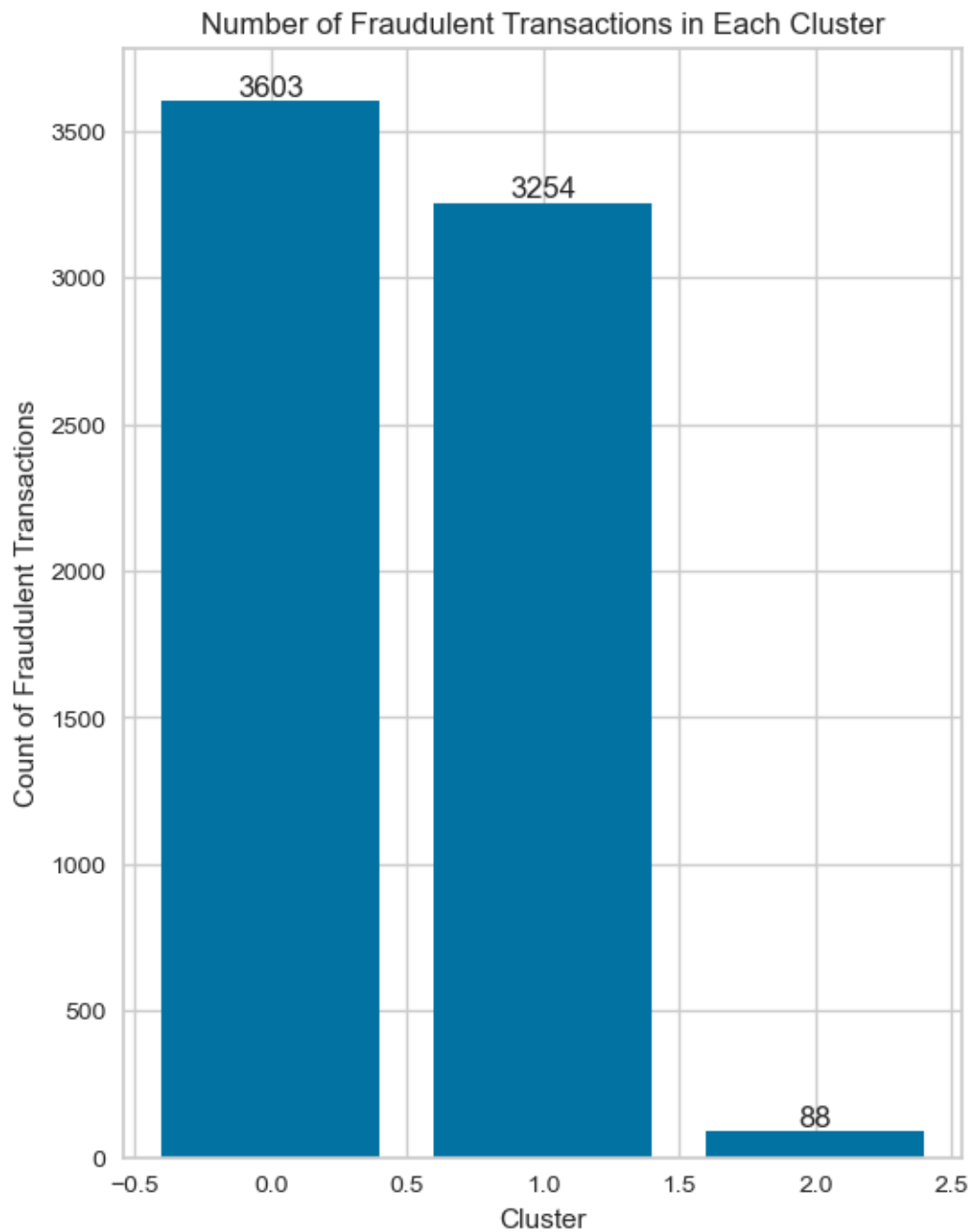


Рисунок 2.25-Кількість шахрайських даних у кожному кластері

Також візуалізуємо кількість даних у кожному з кластерів:

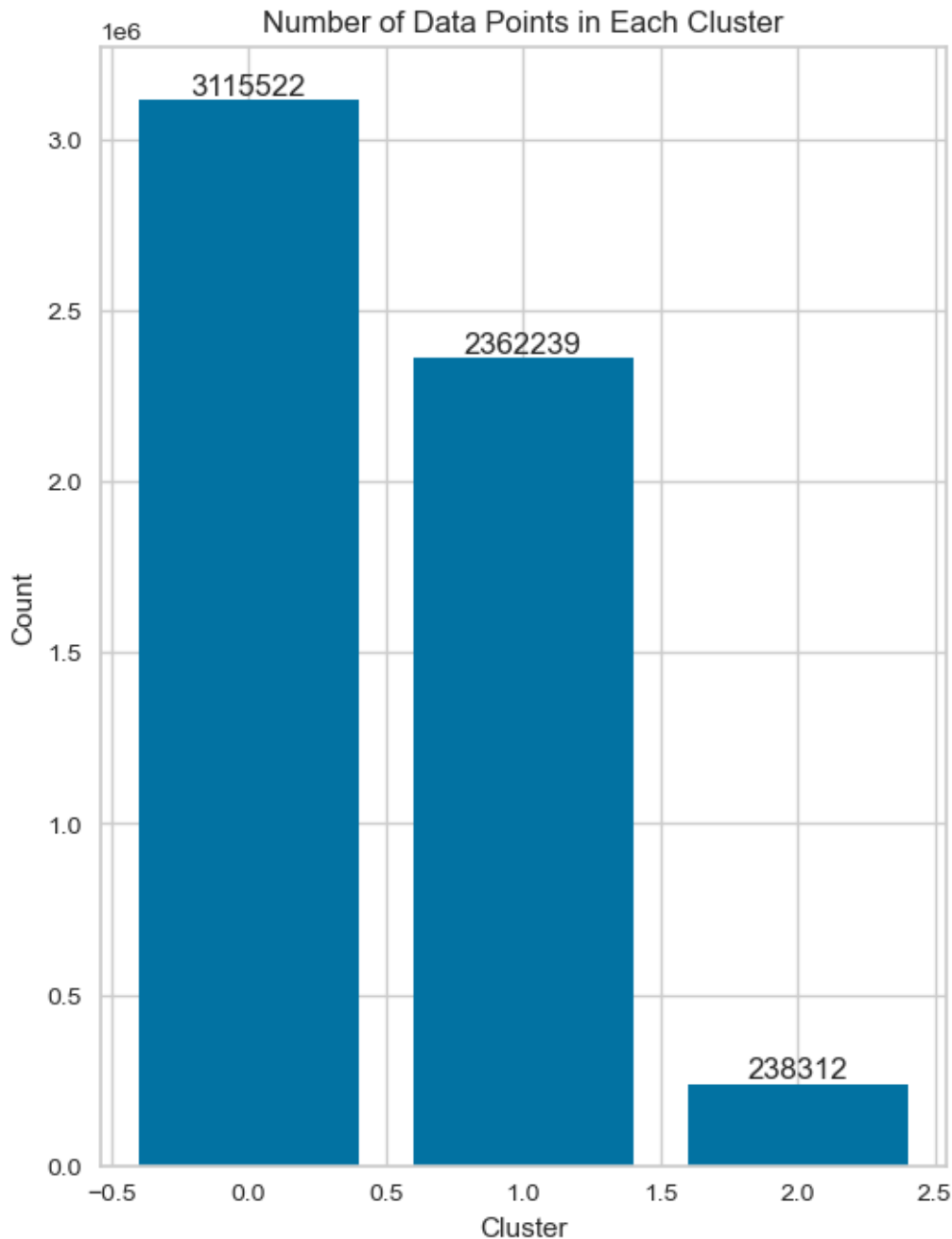


Рисунок 2.26-Кількість даних у кожному кластері

Тепер виконуємо наступні дії:

1. Створюється новий навчальний набір, який включає лише точки даних, що належать до кластера 0. Відповідні мітки класів також відображаються у відповідній змінній.
2. Отримуються кластерні мітки для кожної точки даних у тестовому наборі.

3. Створюється новий тестовий набір, який включає лише точки даних, що належать до кластера 0. Відповідні мітки класів також відображаються у відповідній змінній.

Отримані набори `X_train_cluster0`, `y_train_cluster0`, `X_test_cluster0` і `y_test_cluster0` містять тільки дані, які відповідають кластеру 0.

Тепер можемо побудувати модель прогнозування на основі даних кластеру №0. Виконуємо наступні дії:

1. Створюється об'єкт класифікатора XGBoost (`xgb_model_clust_0`).

2. Проводиться навчання моделі на навчальних даних з використанням валідаційного набору, що вказаний (`eval_set=[(X_test_cluster0, y_test_cluster0)]`). `early_stopping_rounds` вказує на кількість ітерацій, протягом яких буде очікуватися покращення метрики якості на валідаційному наборі перед зупинкою навчання. `verbose=1` вказує на виведення інформації про процес навчання.

3. Здійснюються передбачення моделі на тестовому наборі, що відповідає кластеру 0 (`y_pred = xgb_model_clust_0.predict(X_test_cluster0)`).

4. Обчислюється точність моделі (`accuracy_score`), логлосс (`log_loss`), матриця плутанини (`confusion_matrix`), звіт про класифікацію (`classification_report`), F1-оцінка (`f1_score`) і AUC-ROC (`roc_auc_score`).

```
[97] validation_0-logloss:0.00110
[98] validation_0-logloss:0.00110
[99] validation_0-logloss:0.00110
Accuracy: 0.9996113295982587
Logloss: 0.0010964476809411875
[[346889  15]
 [ 120  314]]
      precision    recall  f1-score   support

     0       1.00      1.00      1.00     346904
     1       0.95      0.72      0.82       434

 accuracy          1.00     347338
 macro avg         0.98     0.86     0.91     347338
weighted avg         1.00     1.00     1.00     347338

F1-score: 0.8230668414154653
AUC-ROC: 0.8617295322596055
```

Рисунок 2.27-Результати моделі прогнозування на основі даних кластеру №0

Результати оцінки моделі XGBoost на тестових даних, що відповідають кластеру 0, є наступними:

- Точність (Accuracy): 0.9996113295982587, що означає, що модель правильно класифікує 99.96% прикладів.
- Логлосс (Logloss): 0.0010964476809411875, який є дуже низьким значенням. Чим менше значення логлосс, тим краще модель прогнозує ймовірності класів.
- Confusion Matrix:
 - True negatives (TN): 346889 - кількість правильно класифікованих негативних прикладів.
 - False positives (FP): 15 - кількість негативних прикладів, які були неправильно класифіковані як позитивні.
 - False negatives (FN): 120 - кількість позитивних прикладів, які були неправильно класифіковані як негативні.
 - True positives (TP): 314 - кількість правильно класифікованих позитивних прикладів.
- F1-оцінка становить 0.8230668414154653 для позитивного класу.
- AUC-ROC: 0.8617295322596055 вказує на те, що модель XGBoost має хорошу дискримінаційну здатність. Це свідчить про те, що модель має високу ймовірність ранжування вибраного випадкового позитивного екземпляру вище, ніж вибраного випадкового негативного екземпляру при встановленні різних порогових значень.

Загалом, модель показує дуже високу точність та хорошу здатність до виявлення фроду (позитивного класу). Проте, такий F1-score буде найменшим з усіх інших моделей, що вказує на гіршу якість цієї моделі.

Далі побудуємо дві нові моделі, на основі лише тренувальних даних та на даних, де кластери є додатковою ознакою у тренувальних та тестових датасетах для порівняння якості моделей.

Модель на основі тренувальних та тестових даних:

```

Accuracy: 0.9997732711928454
Logloss: 0.0007949820151272467
[[634325  23]
 [  121  651]]
      precision    recall  f1-score   support

     0         1.00      1.00      1.00    634348
     1         0.97      0.84      0.90      772

   accuracy          1.00    635120
  macro avg          0.98    635120
 weighted avg          1.00    635120

F1-score: 0.900414937759336
AUC-ROC: 0.9216139955019214

```

Рисунок 2.28-Результат прогнозування на навчальних та тестових даних

- Точність (Accuracy): 0.9997732711928454, що означає, що модель правильно класифікує 99.98% прикладів.
- Логлосс (Logloss): 0.0007949820151272467, який є дуже низьким значенням.
- Confusion Matrix:
 - True negatives (TN): 634325 - кількість правильно класифікованих негативних прикладів.
 - False positives (FP): 23 - кількість негативних прикладів, які були неправильно класифіковані як позитивні.
 - False negatives (FN): 121- кількість позитивних прикладів, які були неправильно класифіковані як негативні.
 - True positives (TP): 651 - кількість правильно класифікованих позитивних прикладів.
 - F1-оцінка становить 0.9004 для позитивного класу.
 - AUC-ROC: 0.9216 вказує на те, що модель XGBoost має хорошу дискримінаційну здатність.

Можемо відразу побачити, що F1-score значно більший за моделі кластера №0.

Тепер побудуємо модель з номерами кластерів у вигляді додаткової ознаки:

```

[97] validation_0-logloss:0.00078
[98] validation_0-logloss:0.00077
[99] validation_0-logloss:0.00077
Accuracy: 0.9997748456984507
[[634327  21]
 [ 122  650]]

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	634348
1	0.97	0.84	0.90	772
accuracy			1.00	635120
macro avg	0.98	0.92	0.95	635120
weighted avg	1.00	1.00	1.00	635120

```

F1-score: 0.9009009009009008
AUC-ROC: 0.9209679035299136

```

Рисунок 2.29-Результат прогнозування з номерами кластерів

Можемо відразу побачити, що f1-score став краще моделі прогнозування без кластеризації на 0.005%, що є незначним покращенням якості, але найточнішу різницю в якості покаже кросс-валідація.

А тепер, побачивши результати моделей XGBoost з різними даними у вигляді кластера № 0 та міток кластерів, побудуємо модель з комбінованими даними кластерів №0 та № 1. Для цього виконуємо фільтрацію навчального та тестового наборів даних з урахуванням кластерних міток. Фільтрація проводиться для кластерів з мітками 0 і 1, при цьому зберігаються всі стовпці ознак.

1. Для навчального набору даних:

- X_train_cluster01 створюється шляхом вибору лише тих рядків з X_train_scaled, де кластерні мітки дорівнюють 0 або 1. Це виконується за допомогою функції np.logical_or, яка створює булевий масив, що містить True, якщо кластерна мітка дорівнює 0 або 1, та False в протилежному випадку.
- y_train_cluster01 містить відповідні категоріальні мітки з y_train_encoded для відфільтрованих рядків у X_train_cluster01. Використовується та сама логіка фільтрації з np.logical_or, щоб вибрати відповідні мітки.

2. Для тестового набору даних:

- `X_test_cluster01` створюється шляхом вибору лише тих рядків з `X_test_scaled`, де кластерні мітки дорівнюють 0 або 1. Аналогічно до навчального набору даних, використовується `np.logical_or` для фільтрації рядків.
- `y_test_cluster01` містить відповідні категоріальні мітки з `y_test_encoded` для відфільтрованих рядків у `X_test_cluster01`. Знову ж таки, використовується та сама логіка фільтрації з `np.logical_or`, щоб вибрати відповідні мітки.

Отже, після виконання цього коду у нас будуть збережені відфільтровані навчальний і тестовий набори даних, які містять лише рядки, що належать до кластерів з мітками 0 або 1, разом з відповідними категоріальними мітками.

Тепер побудуємо модель прогнозування XGBoost на цих даних:

```
[37] validation_0-logloss:0.00189
[38] validation_0-logloss:0.00186
[39] validation_0-logloss:0.00163
Accuracy: 0.9996449762495685
Logloss: 0.0016119901138699122
[[607610  39]
 [ 177  584]]
      precision    recall  f1-score   support

     0         1.00      1.00      1.00     607649
     1         0.94      0.77      0.84         761

 accuracy          1.00          1.00          1.00     608410
 macro avg          0.97          0.88          0.92     608410
 weighted avg          1.00          1.00          1.00     608410

F1-score: 0.8439306358381503
AUC-ROC: 0.8836735595653422
```

Рисунок 2.30-Модель прогнозування на даних з кластера №0 та №1

Бачимо, що маємо один з гірших результатів усіх інших моделей, причиною цьому може бути перенавчання моделі. Можемо змінити деякі параметри моделі, щоб подивитись, чи можемо ми покращити цей результат:

```

[430] validation_0-logloss:0.00068
[431] validation_0-logloss:0.00068
[432] validation_0-logloss:0.00068
Accuracy: 0.9997731792705576
Logloss: 0.0006820731487486438
[[607632 17]
 [ 121 640]]
      precision    recall  f1-score   support

     0       1.00      1.00      1.00    607649
     1       0.97      0.84      0.90      761

 accuracy
macro avg      0.99      0.92      0.95    608410
weighted avg   1.00      1.00      1.00    608410

F1-score: 0.9026798307475318
AUC-ROC: 0.9204853546311139

```

Рисунок 2.31-Глибина 5, кількість дерев – 500, learning_rate = 0.1, subsample =1, colsample_bytree = 1, tree_method = 'auto'

Можемо побачити, що зміна глибини з стандартних 3 до 5, та кількість дерев з 100 до 500 покращили якість моделі.

Також для порівняння самої моделі XGBoost, зробимо модель прогнозування Random Forest на тренувальних та тестових даних:

```

Accuracy: 0.999685098878952
Logloss: 0.03819711016605843
[[634336 12]
 [ 188 584]]
      precision    recall  f1-score   support

     0       1.00      1.00      1.00    634348
     1       0.98      0.76      0.85      772

 accuracy
macro avg      0.99      0.88      0.93    635120
weighted avg   1.00      1.00      1.00    635120

F1-score: 0.8538011695906432
AUC-ROC: 0.8782288834382632

```

Рисунок 2.32-Результат Random Forest

Результат цієї моделі показує один з гірших результатів. Але на цьому прикладі можемо побачити, що результат точності у всіх моделей майже 99,99%, що свідчить про те, що точність не є значимою мірою якості моделі.

Тепер розглянемо, які ознаки були найбільш значущі у процесі прогнозування моделей:

Кластер №0:

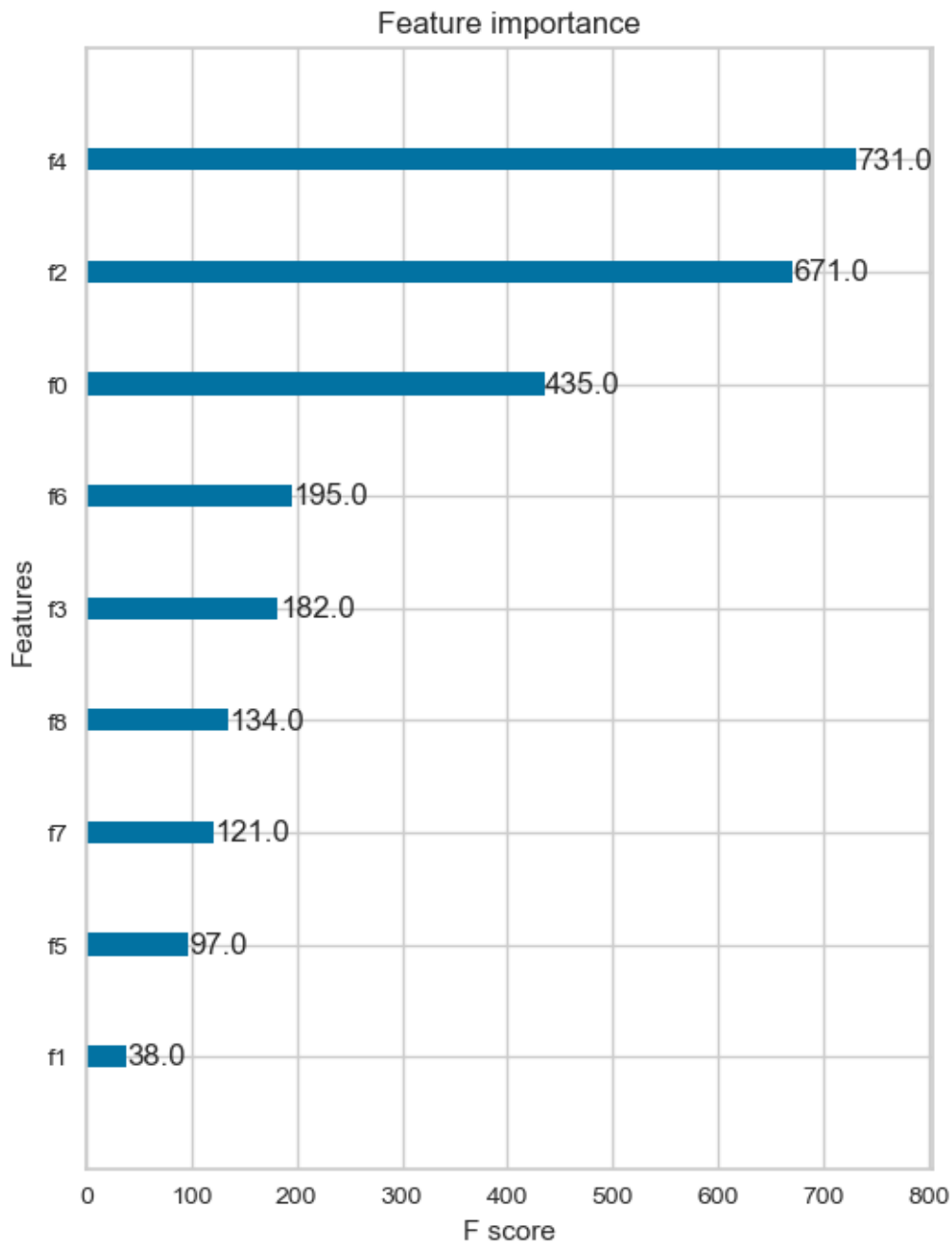


Рисунок 2.33-Feature Importance моделі кластера №0

Бачимо, що найбільш значущими ознаки були f2 (amount) та f4 (oldbalanceOrig), що свідчить про сильний зв'язок цих ознак з шахрайствами, що буде корисною інформацією для рішення задач про виявлення фроду.

Модель на даних тренувального та тестового датасетах:

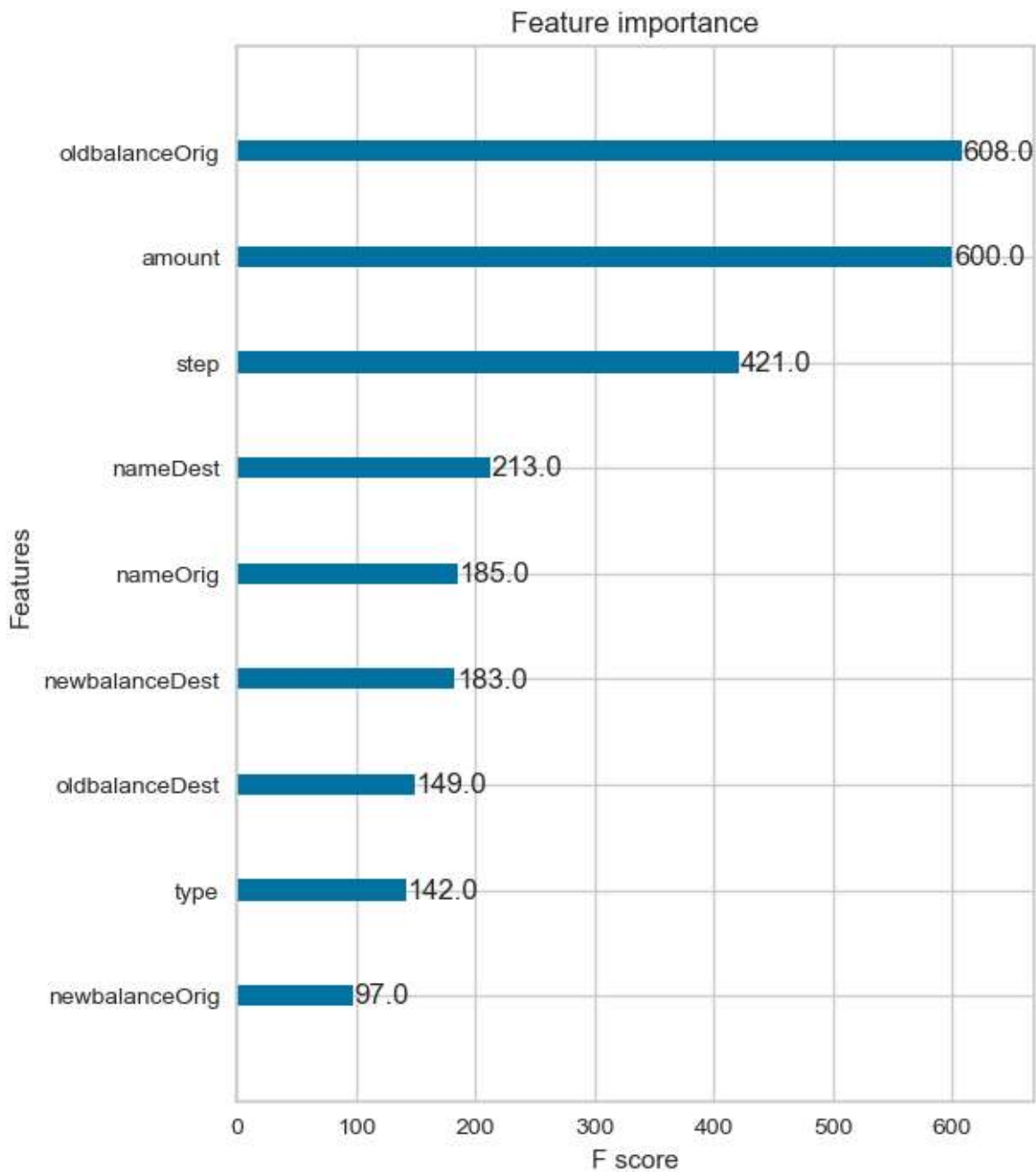


Рисунок 2.34-Feature Importance моделі на тренувальних та тестових даних

Бачимо майже ідентичні зв'язки, але в цьому методі `oldbalanceOrig` менш сильний зв'язок.

Назви кластерів у вигляді додаткової ознаки:

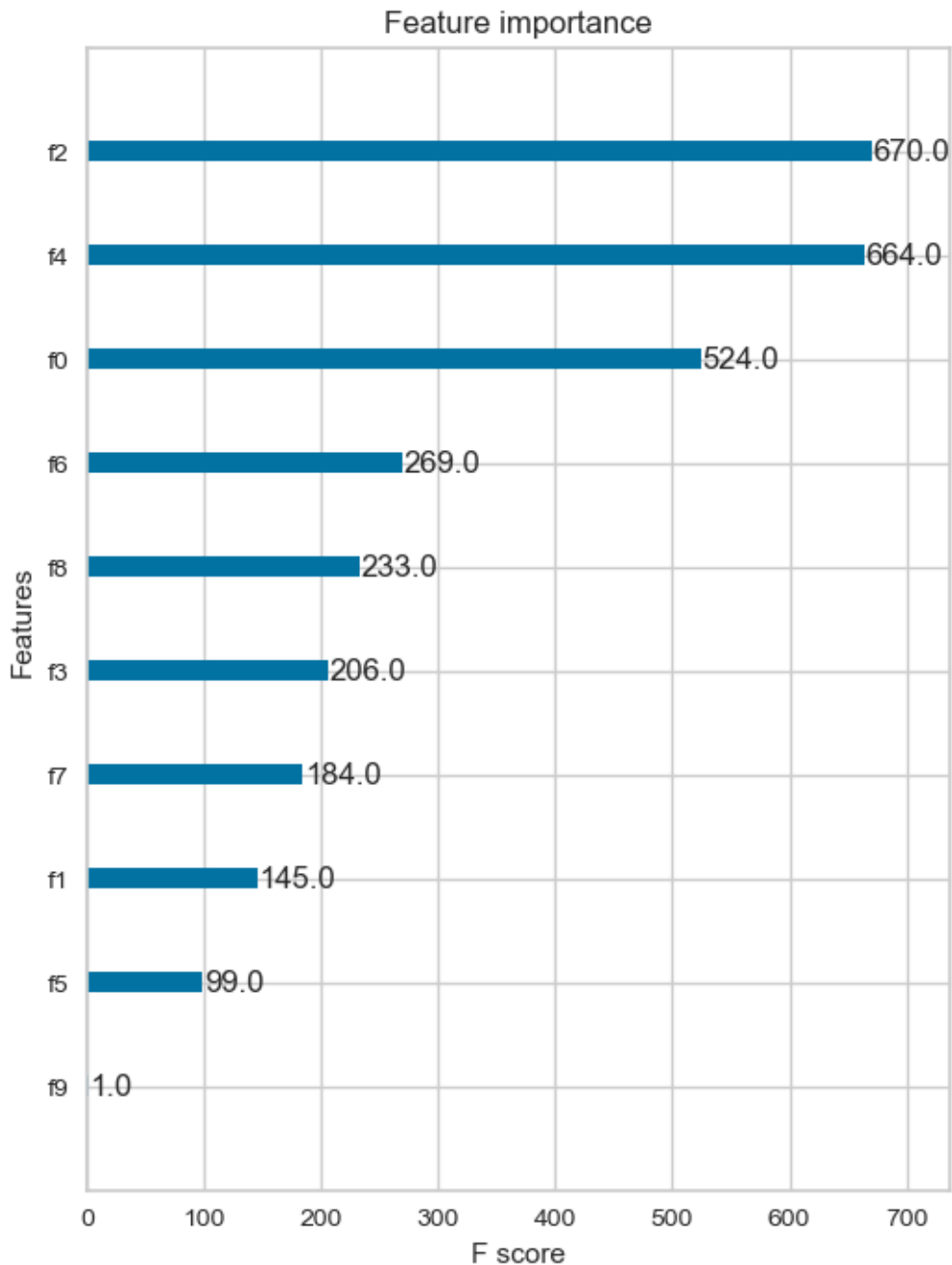


Рисунок 2.35-Feature Importance моделі з додатковою ознакою у вигляді назв кластерів.

Тут можемо побачити, що кластери не вплинули на результат моделі.

Виконаємо кросс-валідацію для моделей оригінальних наборів даних, модель К-середніх, модель комбінованих даних та random Forest та порівняємо їх за допомогою візуалізації.

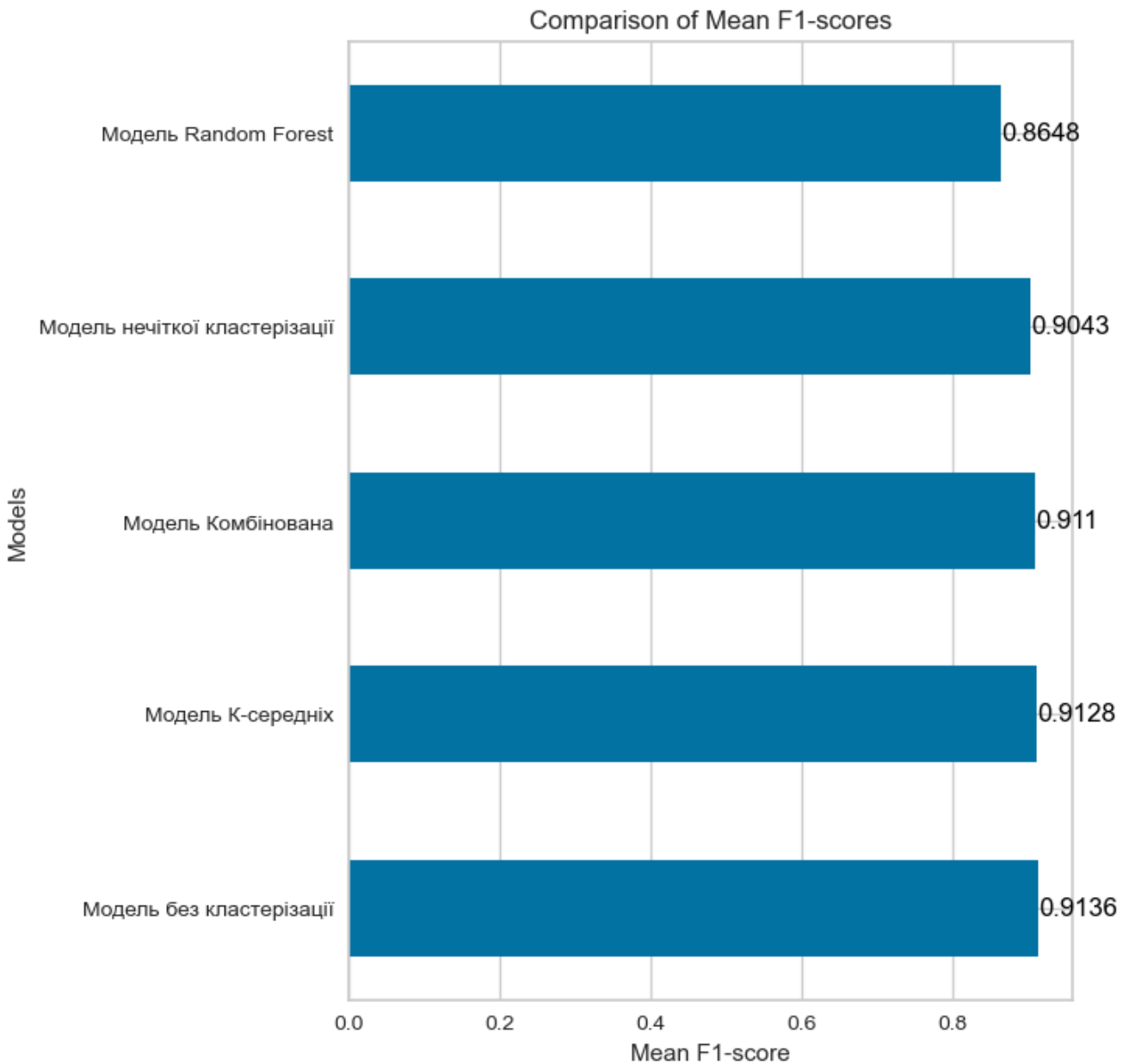


Рисунок 2.36-Візуалізація кросс-валідації моделей прогнозування

Кросс-валідація показує, що найкращою моделлю якості є модель без кластеризації, але їх різниця є дуже малою, якщо порівнювати з моделлю комбінованою та К-середніх, а також гіршу якість Random Forest.

Але це є результатами середнього значення між результатами кросс-валідації, тому подивимося на масив деяких моделей прогнозування, а саме модель без кластеризації та комбіновану:

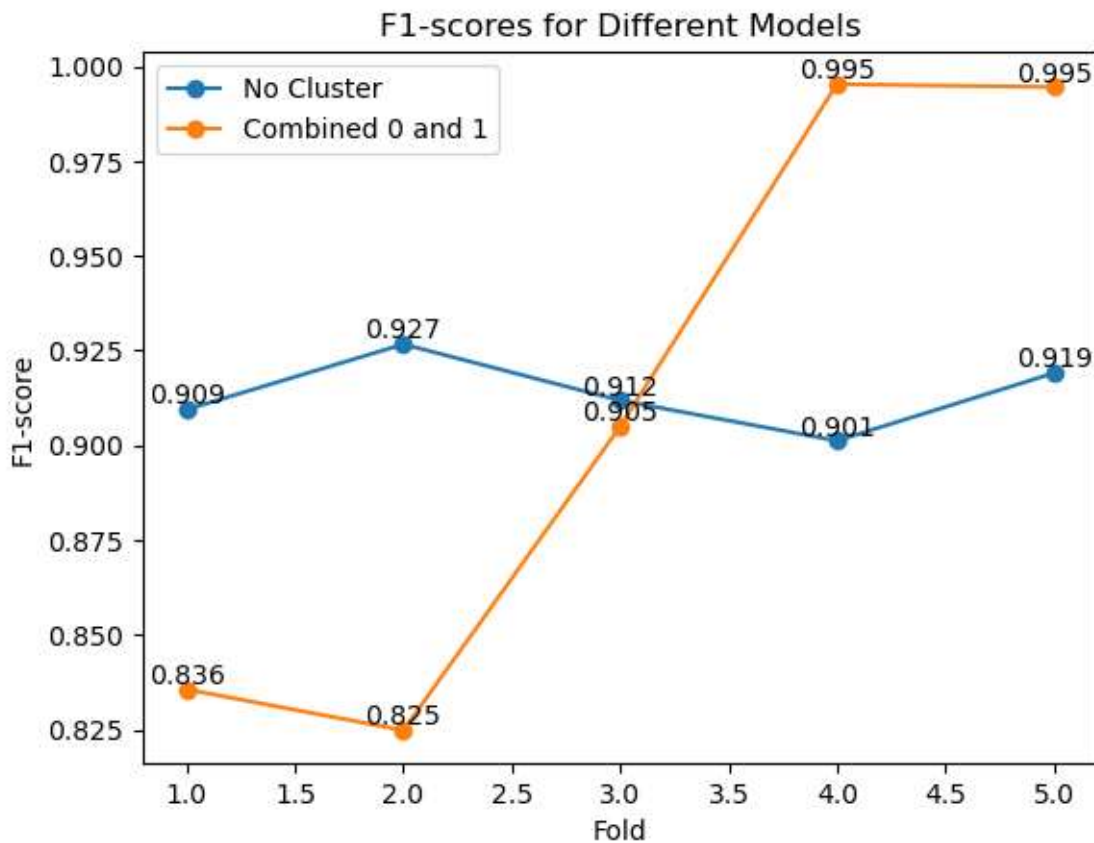


Рисунок 2.37-Масиви кросс-фалідацій для моделей без кластеризації та комбінованої

Ці оцінки представляють ефективність моделі для кожного окремого тестового набору даних. Вони вказують на те, наскільки добре модель працює для кожного конкретного випадку під час кросс-валідації. Більші значення F1-оцінки вказують на кращу ефективність моделі при класифікації.

Бачимо, що комбінований метод дає у останніх тестах найбільший результат, що свідчить про вдосконалення якості моделі при подальших тестах, але при узагальненні, кращою стає модель без кластеризації за рахунок збалансованості її масиву.

Далі розглянемо результати нечіткої кластеризації.

2.4.2 Метод Fuzzy Clustering

Спочатку обчислимо оптимальну кількість кластерів для кластеризації, для цього виконаємо наступні дії:

1. Створюється вибірка даних за допомогою функції `make_blobs`. Ця функція генерує зразки даних у формі кластерів, де кожен кластер представляється як гаусіанський розподіл.

2. Визначається діапазон значень k , які будуть використовуватись для тестування алгоритму кластеризації. У даному випадку, діапазон складається з чисел від 2 до 9.

3. Створюється пустий список `f_indices`, який буде містити значення F-індексу для кожного значення k .

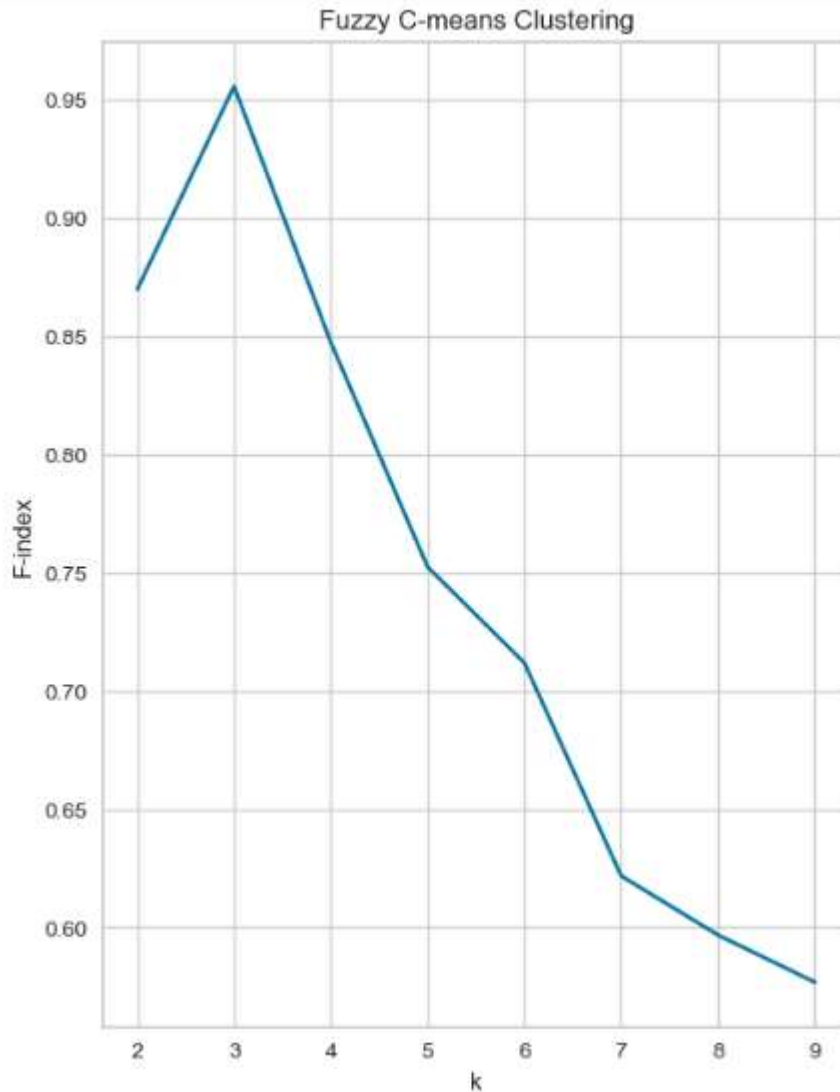
4. Виконується цикл, в якому обчислюється F-індекс для кожного значення k . Для кожного k , використовуючи алгоритм Fuzzy C-means, обчислюються центри кластерів `cntr`, матриця приналежності `u`, а також значення F-індексу `f`. Значення F-індексу додається до списку `f_indices`.

5. Побудова графіка залежності F-індексу від значення k . На графіку вісь X представляє значення k , а вісь Y - значення F-індексу.

6. Визначення оптимального значення k як індексу з найбільшим значенням F-індексу. Оскільки діапазон k починається з 2, то до індексу додається 2.

7. Виведення на екран оптимального значення k .

Цей код виконує пошук оптимального значення k для алгоритму Fuzzy C-means, використовуючи F-індекс як метрику оцінки якості кластеризації. Це дозволяє визначити кількість оптимальних кластерів в даних.



The optimal number of clusters is 3

Рисунок 2.38-Оптимальна кількість кластерів

Далі проведемо наступні операції:

1. Імпорт бібліотеки `LabelEncoder` з модулю `sklearn.preprocessing`. `LabelEncoder` використовується для перетворення категоріальних міток в числовий формат.
2. Створення екземпляра `LabelEncoder` - об'єкту, який виконує перетворення.
3. Застосування методу `fit_transform` до категоріальних міток у `y_train`. Метод `fit_transform` спочатку навчає `LabelEncoder` на унікальних значеннях міток у `y_train` і потім застосовує перетворення до `y_train`, замінюючи категоріальні мітки числовими значеннями.

4. Застосування методу `transform` до категоріальних міток у `y_test`. Метод `transform` застосовує перетворення, навчене на `y_train`, до `y_test`, замінюючи категоріальні мітки числовими значеннями.

5. Імпорт бібліотеки `skfuzzy` з модулю `skfuzzy`. Бібліотека `skfuzzy` використовується для роботи з нечіткими множинами та нечітким кластеруванням.

6. Встановлення кількості кластерів `n_clusters` на значення 3.

7. Застосування алгоритму нечіткого кластерування C-середніх (`fuzzy C-means`) до `X_train`. Метод `fuzz.cluster.cmeans` приймає матрицю `X_train` (транспоновану) та інші параметри, такі як кількість кластерів, максимальна кількість ітерацій та допустима похибка. Результати включають центри кластерів (`cntr_train`) та матрицю приналежності (`u_train`), що містить ймовірності приналежності кожного зразка до кожного кластера.

8. Застосування отриманих центрів кластерів (`cntr_train`) до `X_test`. Метод `fuzz.cluster.cmeans_predict` використовує отримані центри кластерів для передбачення приналежності зразків з `X_test` до кластерів. Результати включають матрицю приналежності (`u_test`), що містить ймовірності приналежності кожного зразка з `X_test` до кожного кластера.

Тепер перевіримо розмірність `u_train` та `u_test`:

```
u_train shape: (5716073,)
u_test shape: (3, 635120)
```

Рисунок 2.39-Розмірність `u_train` та `u_test`

За допомогою функції `predict` виведемо масив `u_test`:

```

(array([[0.02140822, 0.02990875, 0.03060859, ..., 0.02917758, 0.02068882,
        0.53564924],
       [0.06341336, 0.10242862, 0.08711466, ..., 0.08478136, 0.06107368,
        0.23210779],
       [0.91517839, 0.86766263, 0.88227675, ..., 0.88604106, 0.9182375 ,
        0.23224297]]) , array([[0.38142722, 0.4541 , 0.39772568, ..., 0.54250066, 0.22213158,
        0.31464057],
       [0.20076836, 0.42265229, 0.01038959, ..., 0.33341433, 0.43322108,
        0.29114175],
       [0.41780441, 0.12324771, 0.59188473, ..., 0.12408502, 0.34464734,
        0.39421768]]) , array([[15846849.65113703, 15745146.88351158, 15771253.47072183, ...,
        15919609.4469782 , 15783719.25028777, 29194920.26679619],
       [ 9207521.06296036,  8508151.1389627 ,  9348506.95798217, ...,
        9339139.39661319,  9186501.66720757, 44350924.94092631],
       [ 2423708.12656927,  2923279.94090551,  2937551.29899943, ...,
        2888884.16009826,  2369188.88206488, 44338015.4294705 ]]) , array([[3.63769933e+19, 9.68629106e+18]], 2, 0.7407809824631
61)

```

Рисунок 2.40-Масив u_test

U_test містить 3 рядки, матрицю потрібно транспонувати та відкинути третій рядок, бо в сумі три стовпці дадуть стовпець з одиниць. Причина, чому сума цих стовпців дорівнює 1, полягає в особливості результатів нечіткої кластеризації.

У нечіткій кластеризації кожен приклад може належати до кожного кластера з певною ступінню належності. Кожне значення у стовпцях матриці u_test представляє ступінь належності прикладу до певного кластера. Загальна сума ступенів належності для кожного прикладу має бути рівна 1.

Тому виправимо u_test:

```

u_test_transposed shape: (2,)
array([[0.06341336, 0.02140822],
       [0.10242857, 0.02990876],
       [0.08711462, 0.03060886 ],
       ...,
       [0.08478132, 0.02917759],
       [0.06107365, 0.02068882],
       [0.2321078 , 0.53564922]])

```

Рисунок 2.41-Виправлений u_test

Далі виконаємо об'єднання матриць X_test і u_test_transposed для створення нової матриці X_test_new. Також він об'єднує матриці X_train і u_train для створення нової матриці X_train_new.

У результуючих матрицях X_test_new і X_train_new ми отримуємо розширену версію вхідних даних, де кожний рядок представляє собою поєднання

ознак з вихідної матриці і ступеня належності до кожного кластера. За допомогою функції `np.hstack()` виконується горизонтальне об'єднання матриць, де кожна матриця розміщується вгороджено одна поруч із однією.

В результуючих матрицях `X_test_new` і `X_train_new` кількість стовпців збільшилася на два, оскільки були додані стовпці зі ступенями належності до кожного кластера. Це буде потрібно для моделі прогнозування, яка використовує дані кожних кластерів як нову ознаку.

Після цього можемо розробити моделі прогнозування.

Модель на даних нечіткої кластеризації:

```
[95] validation_0-logloss:0.00077
[96] validation_0-logloss:0.00077
[97] validation_0-logloss:0.00077
[98] validation_0-logloss:0.00077
[99] validation_0-logloss:0.00077
Accuracy: 0.999770122181635
Logloss: 0.0007661827449868437
[[634322  26]
 [ 120   652]]
      precision    recall  f1-score   support

     0         1.00      1.00      1.00     634348
     1         0.96      0.84      0.90       772

 accuracy         1.00         1.00         1.00     635120
 macro avg         0.98         0.92         0.95     635120
 weighted avg         1.00         1.00         1.00     635120

 F1-score: 0.8993103448275863
 AUC-ROC: 0.9222592992630416
```

Рисунок 2.42-Результат моделі прогнозування на даних нечіткої кластеризації

Результати моделі мають майже такі ж самі результати з моделями К-середніх.

Бачимо більш-менш такий ж самий результат в порівнянні з моделлю методу К-середніх.

Маючи результати моделей на даних нечіткої кластеризації, К-середніх, комбінованого та без кластеризації порівняємо результати загальних моделей:

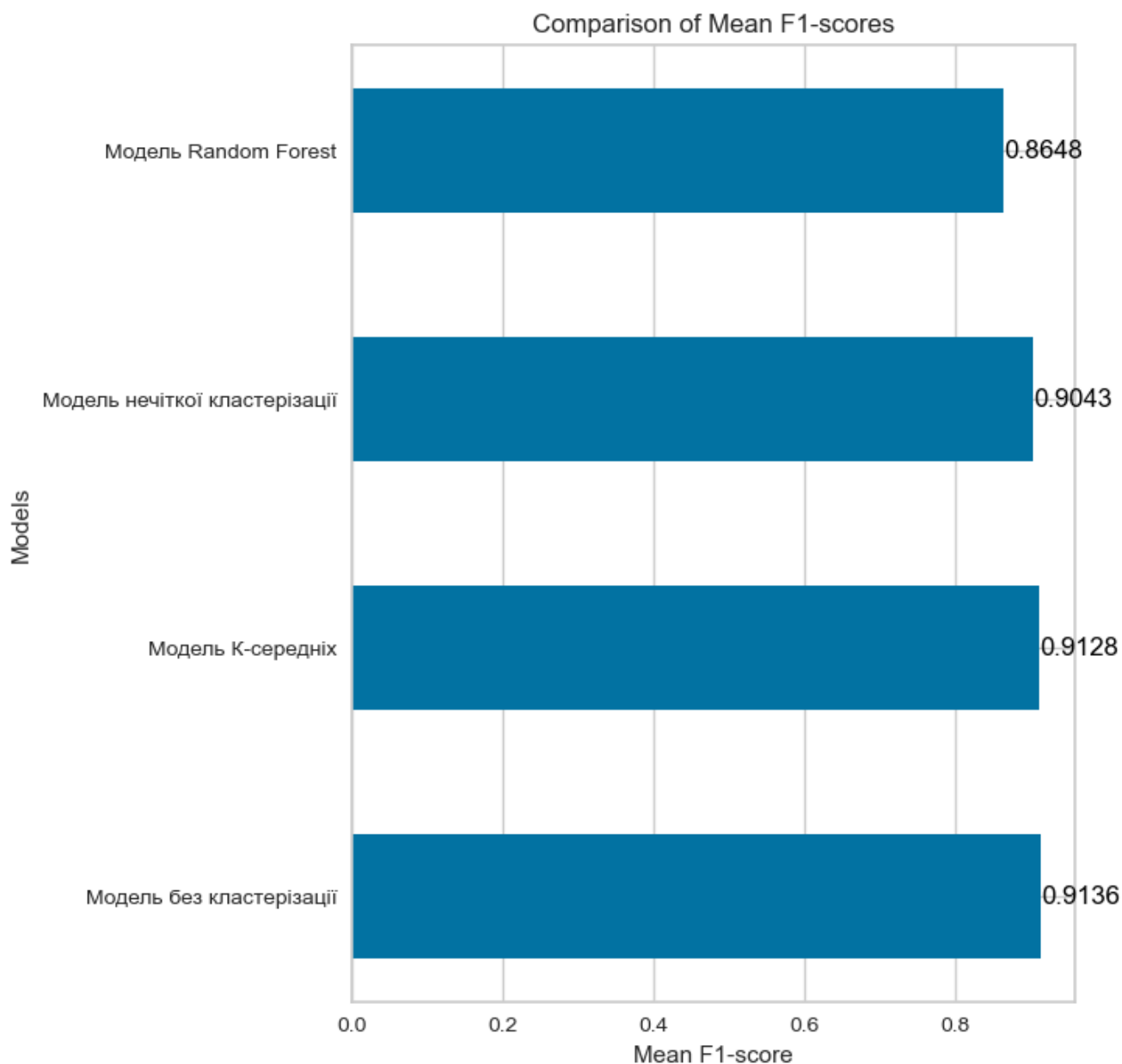


Рисунок 2.43-Результати загальних моделей прогнозування

Бачимо на графіку, що найкращу якість має модель без кластеризації, випереджаючи модель K-середніх на лише 0.0008%, що є незначною, але явною різницею на графіку між цими моделями. Таким чином ми визначили, що модель без кластеризації позитивно впливає на якість моделі прогнозування XGBoost. Та важливо зазначити, що найкращий потенціал до підвищення якості моделі прогнозування у подальших тестах з даними має модель з комбінованими даними кластерів №0 та №1 через те, що вона має вищі значення F1-оцінки в порівнянні з моделлю без використання кластерів.

2.4.3 Метод CatBoost

Загально, метод CatBoost є менш вимогливим до даних датасету, так як він був розроблений для роботи з категоріальними змінними. На відміну від інших алгоритмів машинного навчання, які вимагають перетворення категоріальних змінних у числовий формат за допомогою одноразового кодування або подібних методів, CatBoost може обробляти ці змінні нативно, що значно спрощує процес підготовки даних і підвищує продуктивність моделі.

Тому, завдяки цим перевагам ми можемо сконцентруватися на самій моделі та її параметрах для аналізу результатів.

Для початку, побудуємо стандартну модель CatBoost для нашого датасету:

```
# Initialize CatBoostClassifier with recommended parameters for large
imbalanced datasets
```

```
model_standart = CatBoostClassifier(
    iterations=2000,          # Large number of trees, can adjust based on time
constraints
    learning_rate=0.05,      # Slower learning rate for better convergence
    depth=8,                 # Moderate depth, balances complexity and speed
    l2_leaf_reg=5,          # Regularization to prevent overfitting
    auto_class_weights='Balanced', # Automatically balances class weights
based on data
    eval_metric='AUC',       # AUC is appropriate for imbalanced
classification
    random_seed=42,
    early_stopping_rounds=100, # Stops early if performance doesn't
improve for 100 rounds
    verbose=100              # Prints update every 100 iterations
)
```

```
# Train the model, specifying categorical features
```

```

model_standart.fit(
    X_train, y_train,
    cat_features=cat_features,          # Specify categorical features for optimal
handling
    eval_set=(X_test, y_test),
    early_stopping_rounds=100          # Stop if no improvement for 100 rounds
)

```

```

0:      test: 0.9893314 best: 0.9893314 (0)      total: 11.4s   remaining: 6h 19m 32s
100:    test: 0.9991104 best: 0.9991116 (98)    total: 13m 34s remaining: 4h 15m 8s
200:    test: 0.9991702 best: 0.9992642 (134)  total: 24m 4s  remaining: 3h 35m 26s
Stopped by overfitting detector (100 iterations wait)

bestTest = 0.9992642235
bestIteration = 134

Shrink model to first 135 iterations.
<catboost.core.CatBoostClassifier at 0x2072cddf2c0>

```

```

AUC Score: 0.9993
Confusion Matrix (Threshold=0.5):
[[1254079  14617]
 [      7   1536]]
Classification Report (Adjusted Threshold):

```

	precision	recall	f1-score	support
0	1.00	0.99	0.99	1268696
1	0.10	1.00	0.17	1543
accuracy			0.99	1270239
macro avg	0.55	0.99	0.58	1270239
weighted avg	1.00	0.99	0.99	1270239

Рисунок 2.44-Результат CatBoost з стандартними параметрами

Основні результати:

1. AUC (Area Under Curve) Score: Значення AUC 0.9993 підтверджує, що модель успішно розрізняє класи. Це особливо важливо для проблем із дисбалансом, оскільки показує, що модель добре ідентифікує як справжні, так і шахрайські транзакції на рівні ймовірностей.
2. Confusion Matrix (плутанинна матриця):

- Клас 0 (нормальні транзакції): 1,254,079 справжніх нормальних транзакцій правильно передбачені, а 14,617 помилково позначені як шахрайські.
- Клас 1 (шахрайські транзакції): 1,536 шахрайських транзакцій правильно передбачено, і лише 7 були пропущені.

3. Classification Report:

- Precision для класу 1 (шахрайські транзакції): Дуже низька (0.10), що вказує на значну кількість помилкових спрацювань (тобто помилкових передбачень шахрайства).
- Recall для класу 1: Досить висока (1.00), що показує, що майже всі справжні шахрайські випадки виявляються, однак при значній кількості хибнопозитивних передбачень.
- F1-Score: Низький F1-скор для класу 1 (0.17) підтверджує дисбаланс між precision і recall. Це є наслідком того, що модель орієнтована на максимальне виявлення шахрайства (високий recall), але з жертвою точності (низький precision).

4. F2 Score: F2-скор на рівні 0.17 підтверджує, що модель більше орієнтована на recall, ніж на precision, що може бути корисним у бізнес-контексті, де пропуск шахрайської транзакції більш шкідливий, ніж помилкове її виявлення.

Загалом, модель добре виконує задачу виявлення шахрайства з акцентом на виявлення всіх шахрайських транзакцій, але precision залишається низьким. Це свідчить про необхідність коригування, наприклад, шляхом тонкого налаштування параметрів моделі, щоб зменшити кількість хибнопозитивних передбачень.

Для цього у нас є декілька варіантів, такі як:

SMOTE(Synthetic Minority Oversampling Technique) або його модифікації, такі як Borderline-SMOTE або SMOTEENN, допоможуть підвищити вагу міноритарного класу шляхом генерації синтетичних прикладів.

Андерсемплінг мажоритарного класу може допомогти зменшити його вплив, однак може призвести до втрати інформації. Цей метод краще працює у поєднанні з оверсемплінгом міноритарного класу.

Метод байєсової оптимізації (наприклад, Optuna або Hyperopt) для пошуку найкращих значень гіперпараметрів. Байєсова оптимізація ітеративно покращує параметри на основі минулих результатів, що дозволяє швидше досягти оптимального налаштування.

Нам потрібен метод, який за резонний час приведе нас автоматично найкращі параметри моделі, які зможуть покращити результати щодо мінорного класу. Для цього підходить метод байєсової оптимізації, а саме OPTUNA:

- Optuna використовує байєсову оптимізацію, яка допомагає скоротити час налаштування параметрів, аналізуючи попередні результати ітерацій. На відміну від випадкового пошуку, який досліджує простір гіперпараметрів випадковим чином, Optuna ітеративно підвищує точність, направляючи пошук у більш перспективні області.
- Optuna підтримує параметр `timeout`, що дозволяє встановити часовий ліміт для пошуку оптимальних гіперпараметрів, що корисно при роботі з великими датасетами. Це дозволяє обмежити пошук оптимальних налаштувань, не перевищуючи доступний час або обчислювальні ресурси.
- Optuna може автоматично знаходити найбільш релевантні значення для параметрів моделі. Наприклад, для великого незбалансованого датасету, Optuna може ітеративно налаштувати такі параметри, як `learning_rate`, `depth`, `l2_leaf_reg`, щоб досягти кращої узагальнюючої здатності.

Через ці переваги, ми можемо задати оптимальних час пошуку параметрів, що зможуть задовольнити цілі цієї задачі:

```
import optuna
from catboost import CatBoostClassifier

def objective(trial):
    param = {
```

```

'iterations': 300,
'learning_rate': trial.suggest_loguniform('learning_rate', 0.01, 0.1),
'depth': trial.suggest_int('depth', 6, 10),
'l2_leaf_reg': trial.suggest_loguniform('l2_leaf_reg', 1, 10),
'bagging_temperature': trial.suggest_uniform('bagging_temperature', 0.5,
1.5),

'auto_class_weights': 'Balanced',
'eval_metric': 'AUC',
'random_seed': 42,
'verbose': 0
}

model = CatBoostClassifier(**param)
model.fit(X_train, y_train, eval_set=(X_test, y_test),
early_stopping_rounds=50, cat_features=cat_features, verbose=0)

auc = model.best_score_['validation']['AUC']

return auc

study = optuna.create_study(direction='maximize')
study.optimize(objective, n_trials=20, timeout=7200) # Limit to 2 hours

print("Best parameters:", study.best_params)
print("Best AUC score:", study.best_value)

```

```

[I 2024-11-04 15:04:27,591] A new study created in memory with name: no-name-c4197266-a535-42e5-b002-464d86e9a692
[I 2024-11-04 15:24:17,914] Trial 0 finished with value: 0.9994291279204909 and parameters: {'learning_rate': 0.062004792956169535, 'depth': 9, 'l2_leaf_reg': 4.875775904244345, 'bagging_temperature': 0.5190726577293099}. Best is trial 0 with value: 0.9994291279204909.
[I 2024-11-04 16:12:46,244] Trial 1 finished with value: 0.9992887878148592 and parameters: {'learning_rate': 0.025713403956418707, 'depth': 9, 'l2_leaf_reg': 5.629696757527235, 'bagging_temperature': 0.781656529214773}. Best is trial 0 with value: 0.9994291279204909.
[I 2024-11-04 17:00:11,560] Trial 2 finished with value: 0.9992797011605436 and parameters: {'learning_rate': 0.01572076472203408, 'depth': 8, 'l2_leaf_reg': 1.914157518800593, 'bagging_temperature': 0.5244624232219238}. Best is trial 0 with value: 0.9994291279204909.
[I 2024-11-04 17:37:56,692] Trial 3 finished with value: 0.9992985833827589 and parameters: {'learning_rate': 0.07235325765853245, 'depth': 10, 'l2_leaf_reg': 9.072066425717649, 'bagging_temperature': 0.91506738700210557}. Best is trial 0 with value: 0.9994291279204909.
Best parameters: {'learning_rate': 0.062004792956169535, 'depth': 9, 'l2_leaf_reg': 4.875775904244345, 'bagging_temperature': 0.5190726577293099}
Best AUC score: 0.9994291279204909

```

Рисунок 2.45-Оптимальні параметри від OPTUNA

Оновлені параметри CatBoost оптимально налаштовують модель для кращого розпізнавання шахрайства в умовах дисбалансу класів:

1. Learning Rate (0.062) – невеликий крок дозволяє моделі точно знаходити оптимальні значення, запобігаючи перенавчанню.

2. Depth (9) – глибина 9 дає змогу моделі враховувати складні залежності, потрібні для розпізнавання шахрайства, без надмірної складності.

3. L2 Leaf Regularization (4.88) – помірна регуляризація запобігає перенавчанню, зберігаючи продуктивність на міноритарному класі.

4. Bagging Temperature (0.52) – обирає різноманітні приклади для кожної ітерації, що покращує узагальнення і збалансованість.

Ці параметри дають змогу моделі ефективно знаходити шахрайські транзакції з високим AUC та recall, мінімізуючи пропущені випадки шахрайства.

Тепер додамо ці параметри до нової моделі та проаналізуємо результати:

```

0:      test: 0.9893314 best: 0.9893314 (0)      total: 13.5s   remaining: 3h 44m 4s
100:    test: 0.9992755 best: 0.9992904 (56)     total: 15m 31s remaining: 2h 18m 10s
200:    test: 0.9994000 best: 0.9994370 (194)    total: 24m 18s remaining: 1h 36m 3s
300:    test: 0.9994240 best: 0.9994519 (290)    total: 37m 39s remaining: 1h 27m 27s
400:    test: 0.9995157 best: 0.9995157 (398)    total: 55m 32s remaining: 1h 22m 58s
500:    test: 0.9995358 best: 0.9995358 (500)    total: 1h 12m 38s remaining: 1h 12m 21s
600:    test: 0.9995454 best: 0.9995526 (565)    total: 1h 26m 33s remaining: 57m 27s
Stopped by overfitting detector. (100 iterations wait)

bestTest = 0.9995525976
bestIteration = 565

Shrink model to first 566 iterations.
): <catboost.core.CatBoostClassifier at 0x20729d6f2c0>

AUC Score: 0.9996
Confusion Matrix:
[[1265794  2902]
 [      13  1530]]
Classification Report:
              precision    recall  f1-score   support

     0             1.00         1.00         1.00    1268696
     1             0.35         0.99         0.51         1543

 accuracy          1.00         1.00         1.00    1270239
 macro avg          0.67         0.99         0.76    1270239
 weighted avg          1.00         1.00         1.00    1270239

```

Рисунок 2.46-Результати оптимізованої моделі

Після налаштування гіперпараметрів моделі CatBoostClassifier на основі результатів Optuna, було отримано покращені результати. Розглянемо основні метрики:

- Показник AUC досяг 0.9996, що вказує на високу здатність моделі розрізняти класи (шахрайські та нешахрайські транзакції) навіть у контексті значного дисбалансу даних.
- Клас 0 (нормальні транзакції): із 1,268,696 прикладів класу 0, модель помилково класифікувала лише 2,902 прикладів (близько 0.2%), що демонструє високу точність для мажоритарного класу.
- Клас 1 (шахрайські транзакції): із 1,543 прикладів класу 1 модель правильно класифікувала 1,530 прикладів, припустившись лише 13 помилок (0.8% від класу 1). Це показує значне покращення в точності для міноритарного класу.

Оновлені параметри моделі значно покращили її здатність розпізнавати шахрайські транзакції з високим recall (0.99) і прийнятним precision (0.35). Це дає змогу виявити майже всі шахрайські операції, що підходить для реального використання, оскільки в такому випадку фінансові ризики будуть мінімізовані завдяки надвисокому recall.

2.4.4 Метод LightGBM

LightGBM часто переважний над XGBoost та CatBoost завдяки високій швидкості навчання, низькому використанню пам'яті та ефективності на великих наборах даних завдяки leaf-wise побудові дерев, histogram-based дискретизації та підтримці категоріальних ознак. Він забезпечує швидшу продуктивність на розподілених системах і дозволяє гнучко налаштовувати гіперпараметри, що робить його ідеальним для обробки дисбалансованих та складних даних з великою кількістю ознак.

Проаналізував поверхневі переваги LightGBM, побудуємо та запустимо модель:

```
# Define parameters for LightGBM model
params = {
    'objective': 'binary',
    'boosting_type': 'gbdt',
    'is_unbalance': True,
    'learning_rate': 0.01, # Lower learning rate
    'num_leaves': 128,    # Increase num_leaves to add model complexity
    'max_depth': 10,     # Increase depth if necessary
    # other parameters remain the same
}

# Convert data to LightGBM dataset format
train_data = lgb.Dataset(X_train, label=y_train)
valid_data = lgb.Dataset(X_test, label=y_test, reference=train_data)

# Train the model with early stopping
model_lightgbm = lgb.train(
    params,
    train_data,
    valid_sets=[valid_data],
    num_boost_round=5000, # Increase this value
    callbacks=[lgb.early_stopping(stopping_rounds=200, verbose=True)]
)

# Predict probabilities and evaluate
y_pred_proba = model_lightgbm.predict(X_test,
num_iteration=model_lightgbm.best_iteration)

# Adjust threshold and evaluate
```



```

threshold = 0.5
y_pred_adjusted = (y_pred_proba >= threshold).astype(int)
print(f"AUC Score: {roc_auc_score(y_test, y_pred_proba):.4f}")
print(f"Confusion Matrix (Threshold={threshold}):\n{confusion_matrix(y_test,
y_pred_adjusted)}")
print("Classification Report (Adjusted Threshold):")
print(classification_report(y_test, y_pred_adjusted))

```

```

Training until validation scores don't improve for 200 rounds
Did not meet early stopping. Best iteration is:
[4999] valid_0's binary_logloss: 0.00104596
AUC Score: 0.9986
Confusion Matrix (Threshold=0.5):
[[1268423    272]
 [   94    1450]]
Classification Report (Adjusted Threshold):

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	1268695
1	0.84	0.94	0.89	1544
accuracy			1.00	1270239
macro avg	0.92	0.97	0.94	1270239
weighted avg	1.00	1.00	1.00	1270239

Рисунок 2.47-Результати LightGBM

- Результати моделі LightGBM для задачі класифікації виглядають дуже високими за різними метриками, що свідчить про хорошу здатність моделі ідентифікувати транзакції як шахрайські (клас 1) або нормальні (клас 0).
- AUC Score становить 0.9986, що є дуже високим показником, і свідчить про здатність моделі добре розрізняти класи навіть при різних значеннях порогів.
- Клас 0 (нормальні транзакції): З 1,268,695 прикладів класу 0 модель помилково класифікувала лише 272 транзакції як шахрайські, що складає лише близько 0.02%. Це свідчить про дуже високу точність для мажоритарного класу, де модель практично не допускає помилок.

- Клас 1 (шахрайські транзакції): З 1,544 прикладів класу 1 модель правильно класифікувала 1,450 прикладів, припустившись лише 94 помилок, або близько 6% від класу 1. Такий показник вказує на високу здатність моделі виявляти шахрайські операції, хоча деякі транзакції все ж залишаються непоміченими, що можна ще покращити для меншості класу.

Високий показник True Positives порівняно з низькою кількістю False Negatives і False Positives свідчить про хорошу здатність моделі виявляти шахрайські транзакції з мінімальними помилками.

Precision класу 1 (шахрайство): 0.84 — модель ідентифікує більшість шахрайських транзакцій точно, але є кілька помилкових спрацьовувань (False Positives).

Recall класу 1 (шахрайство): 0.94 — модель має високу здатність виявляти шахрайські випадки, пропускаючи лише 6% шахрайських транзакцій.

F1-Score класу 1: 0.89 — високий показник балансує precision і recall, що демонструє загальну точність для меншості класу.

Precision та Recall для класу 0 (не-шахрайство) майже ідеальні, що показує, що модель дуже добре розпізнає нормальні транзакції.

LightGBM показав високу точність і хорошу чутливість до шахрайських транзакцій, однак він дещо відстає від XGBoost і CatBoost за деякими метриками, такими як precision та recall для класу 1. Однак, LightGBM все ще є сильною моделлю і може бути ефективним вибором у різних умовах, залежно від важливості конкретних метрик для задачі.

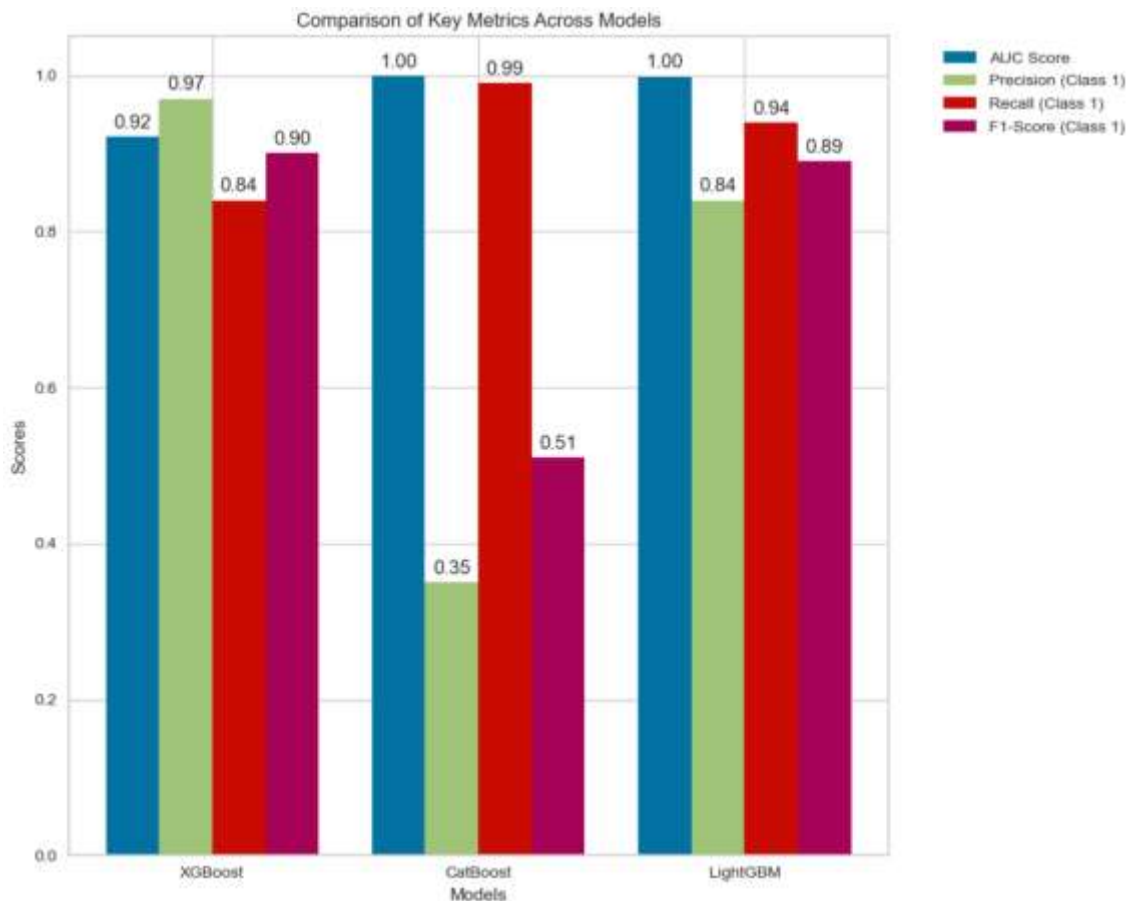


Рисунок 2.48-Порівняння результатів

CatBoost та XGBoost продемонстрували кращі результати в порівнянні з LightGBM у задачі виявлення шахрайства з кількох основних причин:

1. AUC (Площа під кривою ROC):

- CatBoost показує найвищий результат AUC (0.9996), що свідчить про його відмінну здатність до класифікації і точне розрізнення позитивних і негативних класів.
- LightGBM (AUC 0.9986) також демонструє хорошу здатність класифікації, тоді як XGBoost має найнижчий AUC (0.9216), що вказує на меншу ефективність у розрізненні класів.

2. Precision (Точність для класу 1):

- XGBoost має найвищий precision (0.97), що означає високу точність у передбаченні позитивних випадків і мінімальну кількість хибних позитивних результатів.

- LightGBM також має хороший precision (0.84), але він все ж поступається XGBoost.
- CatBoost має найнижчий precision (0.35), що свідчить про більшу кількість хибних позитивних результатів у прогнозах цієї моделі.

3. Recall (Чутливість для класу 1):

- CatBoost демонструє найвищий recall (0.99), що означає здатність дуже точно виявляти позитивні випадки, мінімізуючи пропущені.
- LightGBM має високий recall (0.94), але XGBoost показує найменший recall (0.84), що свідчить про більше пропущених позитивних випадків у його прогнозах.

4. F1-Score (Збалансоване середнє між precision і recall):

- XGBoost має найкращий F1-Score (0.90), що вказує на хороший баланс між точністю і чутливістю.
- LightGBM також має гарний F1-Score (0.89), але CatBoost показує значно нижчий результат (0.51), що свідчить про слабкий баланс між precision та recall.

Таким чином, CatBoost є найкращим варіантом для задач, де важливо досягти максимальної чутливості (recall), оскільки модель має найвищий показник recall і дуже високу AUC. Однак її низький precision може призвести до великої кількості хибних позитивних результатів. Також треба зауважити, що XGBoost забезпечує найкращі результати за precision та F1-Score, що робить її ідеальним вибором для ситуацій, де важливо досягти балансу між точністю та чутливістю, з мінімізацією хибних позитивних результатів. Наостанок, LightGBM також демонструє високі результати по precision та recall, що робить її хорошим варіантом для завдань, де потрібен збалансований підхід.

Загалом, вибір моделі залежить від конкретних вимог: якщо важливий баланс між точністю та чутливістю, то найкращим вибором є XGBoost; для мінімізації пропущених позитивних випадків — CatBoost.

ВИСНОВКИ

У цій кваліфікаційній роботі було проведено дослідження та порівняння трьох методів класифікації — CatBoost, XGBoost та LightGBM — для задачі виявлення шахрайства в фінансових транзакціях. Тестування моделей проводилося на стандартних тренувальних та тестових наборах даних, включаючи незбалансований набір, що є типовим для задач у сфері виявлення шахрайства. Вибір моделей обґрунтовано їх здатністю працювати з такими даними, де одна з категорій (наприклад, шахрайські транзакції) значно рідша за іншу.

XGBoost продемонстрував високі результати за основними метриками. Зокрема, модель досягла значного значення precision для класу 1 (шахрайства), яке склало 0.97. Це вказує на високу точність у передбаченні шахрайських транзакцій, з мінімальною кількістю хибних позитивних результатів. Крім того, XGBoost показав високі показники F1-Score (0.90) та AUC (0.9216), що вказує на здатність моделі зберігати хороший баланс між точністю та чутливістю, а також ефективно розрізняти позитивні і негативні випадки. Найкращі результати були досягнуті завдяки оптимізації гіперпараметрів, що дозволило значно покращити точність для меншості класів, зберігаючи при цьому високу ефективність на більшості даних.

CatBoost також показав дуже хороші результати, зокрема в частині recall, досягнувши показника 0.99, що свідчить про здатність моделі майже без втрат виявляти всі шахрайські транзакції. Також варто відзначити високе значення AUC (0.9996), яке свідчить про хорошу загальну ефективність моделі. Однак precision для класу 1 в CatBoost був значно нижчим (0.35 при порозі 0.5), що вказує на те, що модель схильна до більшої кількості хибних позитивних результатів, хоча при цьому вона добре справляється з виявленням шахрайських транзакцій. Важливо відзначити, що CatBoost ефективно обробляє категоріальні ознаки без потреби в додатковому попередньому обробленні, що є важливою перевагою при роботі з реальними наборами даних.

LightGBM показав досить хороші результати, зокрема високе значення AUC (0.9986) та хороший precision (0.84), що вказує на здатність моделі правильно класифікувати більшість позитивних випадків. Проте, recall LightGBM був трохи нижчим (0.94), що означає, що модель пропускає деякі шахрайські транзакції порівняно з CatBoost та XGBoost. Водночас, binary logloss у LightGBM був дуже низьким, що свідчить про високу точність прогнозування ймовірностей для кожного класу. Загалом, LightGBM є конкурентоспроможною моделлю, проте її результати в порівнянні з іншими двома методами, зокрема за показниками precision та recall, виявилися дещо нижчими.

У підсумку, в рамках цієї роботи було визначено, що найкращі результати для задачі виявлення шахрайства виявляються у XGBoost, завдяки високому precision, гарному F1-Score та хорошому балансі між точністю і чутливістю. CatBoost є найкращим вибором для задач, де важливо максимізувати recall, оскільки модель має найвищий показник чутливості. LightGBM, хоча і демонструє хороші результати в загальному, має дещо нижчі показники в порівнянні з іншими двома методами. Однак усі три моделі є потужними інструментами для розв'язання задач виявлення шахрайства, і їх застосування залежить від конкретних вимог до точності та чутливості в конкретному проекті.

Результати дослідження підтверджують, що підходи, використані в роботі, є ефективними для вирішення задач з незбалансованими даними в сфері виявлення шахрайства. Рекомендовано використовувати XGBoost для задач, де важливий баланс точності та чутливості, а також CatBoost для задач, що вимагають високої чутливості при допустимій кількості хибних позитивних результатів.

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Кваліфікаційна робота магістра [Електронний ресурс] : методичні рекомендації для здобувачів ступеня магістра освітньо-професійної програми «Системний аналіз» зі спеціальності 124 Системний аналіз / уклад.: Т.А. Желдак, Т.В. Хом'як, А.В. Малієнко ; М-во освіти і науки України, Нац. техн. ун-т «Дніпровська політехніка». – Дніпро : НТУ «ДП», 2024. – 33 с. <https://ir.nmu.org.ua/handle/123456789/167921>
2. ADASYN: Adaptive Synthetic Sampling Approach for Imbalanced Learning (https://www.researchgate.net/publication/224330873_ADASYN_Adaptive_Synthetic_Sampling_Approach_for_Imbalanced_Learning).
3. Evaluating XGBoost for Balanced and Imbalanced Data: Application to Fraud Detection (https://www.researchgate.net/publication/369556752_Evaluating_XGBoost_for_Balanced_and_Imbalanced_Data_Application_to_Fraud_Detection).
4. A NEW APPROACH TO CUSTOMER CHURN MODELLING USING LLM AND CATBOOST Ankit Chopra*1, Deependra Singh Rathore*2, Ayush Kumar*3 - (https://www.irjmets.com/uploadedfiles/paper/issue_5_may_2024/57236/final/fin_irjmets1718378031.pdf)
5. A Novel Prediction Model for Diabetes Detection Using Gridsearch and A Voting Classifier between Lightgbm and KNN Nachiket Dunbray Rashmi, Rane Sparsh Nimje, Jayesh Katade, Shreyas Mavale – (https://www.researchgate.net/publication/356081300_A_Novel_Prediction_Model_for_Diabetes_Detection_Using_Gridsearch_and_A_Voting_Classifier_between_Lightgbm_and_KNN)
6. Random Forests (https://www.researchgate.net/profile/Adele-Cutler/publication/236952762_Random_Forests/links/0a85e52f509178323a000000/Random-Forests.pdf?origin=publication_detail).

7. IMBENS: Ensemble Class-imbalanced Learning in Python (https://www.researchgate.net/publication/356602000_IMBENS_Ensemble_Class-imbalanced_Learning_in_Python).

8. Cost-Sensitive Learning vs. Sampling: Which is Best for Handling Unbalanced Classes with Unequal Error Costs? (https://www.researchgate.net/publication/220705031_Cost-Sensitive_Learning_vs_Sampling_Which_is_Best_for_Handling_Unbalanced_Classes_with_Unequal_Error_Costs).

9. Addressing the Curse of Imbalanced Training Sets: One-Sided Selection (<https://sci2s.ugr.es/keel/pdf/algorithm/congreso/kubat97addressing.pdf>).

10. Learning from Imbalanced Data by He, H., & Garcia, E. A. (2009) (https://www.academia.edu/29164932/Learning_from_Imbalanced_Data).

11. Classification of Imbalanced Data: Review of Methods and Applications (https://www.researchgate.net/publication/350084459_Classification_of_Imbalanced_Data_Review_of_Methods_and_Applications).

12. Unstructured Data Analytics: How to Improve Customer Acquisition, Customer Retention, and Fraud Detection and Prevention – by Jean Paul Isson – 26ct.

13. Resampling strategies for imbalanced datasets (<https://www.kaggle.com/rafjaa/resamplingstrategies-for-imbalanced-datasets>).

14. AFS (Advanced Fraud Detection System)– AEC, 2020. – (<https://www.aec.cz/en/afs>).

15. The Only Complete Online Banking Fraud Prevention Solution (<https://www.threatmark.com>).

16. "Mastering Machine Learning with Python in Six Steps: A Practical Implementation Guide to Predictive Data Analytics Using Python" by Manohar Swamynathan.

17. Chen, T., Guestrin, C. (2016). "XGBoost: A Scalable Tree Boosting System." - <https://paperswithcode.com/paper/xgboost-a-scalable-tree-boosting-system>

18. Xu, Y., Zha, L., Li, X. (2019). "Application of XGBoost in Credit Card Fraud Detection." Journal of Data Science and Analysis. - <https://link.springer.com/journal/41060>
19. Dorogush, A. V., Ershov, V., Gulin, A. (2018). "CatBoost: unbiased boosting with categorical features." Advances in Neural Information Processing Systems (NeurIPS) - <https://www.semanticscholar.org/paper/CatBoost%3A-unbiased-boosting-with-categorical-Ostroumova-Gusev/ee0a0f04d45f86bf50b24d7258e884725fcaa621>
20. Prokhorenkova, L., Gusev, G., Vorobev, A., Dorogush, A. V., Gulin, A. (2018). "CatBoost: gradient boosting with categorical features support." IEEE International Conference on Artificial Intelligence - <https://join.if.uinsgd.ac.id/index.php/join/article/view/1324>
21. Практикум з диференційних рівнянь [Електроний ресурс] : навчальний посібник / Л.С. Коряшкіна, О.Д. Станіна, Ю.О. Шевченко; М-во освіти і науки України, Нац. техн. ун-т «Дніпровська політехніка» - Дніпро : НТУ «ДП», 2024 – 178 с. <https://ir.nmu.org.ua/handle/123456789/167658>
22. Аналіз та обробка великих даних [Електронний ресурс] : методичні рекомендації до виконання практичних робіт для здобувачів ступеня магістра освітньо-професійної програми «Системний аналіз» зі спеціальності 124 Системний аналіз / М-во освіти і науки України, Нац. техн. ун-т «Дніпровська політехніка». – Дніпро : НТУ «ДП», 2024. – 82 с. <https://ir.nmu.org.ua/handle/123456789/167968>
23. Машинне навчання [Електронний ресурс] : методичні рекомендації до виконання практичних робіт для здобувачів ступеня магістра освітньо-професійної програми «Системний аналіз» зі спеціальності 124 Системний аналіз / уклад.: Т.А. Желдак, О.Б. Владико, А.В. Малієнко, Д.М. Гаранжа ; М-во освіти і науки України, Нац. техн. ун-т «Дніпровська політехніка». – Дніпро : НТУ «ДП», 2024. – 48 с. <https://ir.nmu.org.ua/handle/123456789/167920>

ДОДАТОК А. Відомість матеріалів кваліфікаційної роботи

№ з/п	Позначення				Найменування	Кількість аркушів	Примітки			
1										
2					Документація					
3										
4	124.КР.24.05.ПЗ				Пояснювальна записка	136	Формат А4			
5										
6	124.КР.24.05.ПЗ				Демонстраційний матеріал	№2	Презентація на CD-R			
7										
8	124.КР.24.05.ПЗ				Копія роботи	1	Диск CD-R			
9										
10										
11										
12										
13										
14										
15										
16										
17										
18										
					124.КР.24.05.ДА.ПЗ.					
Змін.	Аркуш	№ докум.	Підпис	Дата						
Розроб.	Жучков С.І				Матеріали кваліфікаційної роботи	Літ.	Аркуш	Аркушів		
К. розд.										
Керівн.						НТУ «ДП», 12; 124М-23-1				
Н.контр.										
Зав. каф.										