

Міністерство освіти і науки України
Департамент науки і освіти Харківської обласної державної адміністрації
Комунальний заклад «Харківська обласна Мала академія наук
Харківської обласної ради»

Відділення інформаційних технологій
Секція: системи та технології штучного інтелекту.

ДОСЛІДЖЕННЯ МЕТОДІВ ВИЯВЛЕННЯ АНОМАЛІЙ У ФІНАНСОВИХ ТРАНЗАКЦІЯХ

Роботу виконав:

Зіненко Михайло Євгенович, учень 11 класу
Комунального закладу «Харківський академічний
ліцей № 45 Харківської міської ради»

Наукові керівники:

Шаповалова Марія Ігорівна, старший викладач
кафедри математичного моделювання та
інтелектуальних обчислень в інженерії навчально-
наукового інституту комп'ютерного моделювання,
прикладної фізики та математики Національного
технічного університету «Харківський
політехнічний інститут», кандидат технічних наук

Арзубов Микола Олексійович, вчитель
інформатики Комунального закладу «Харківський
академічний ліцей № 45 Харківської міської ради»

ДОСЛІДЖЕННЯ МЕТОДІВ ВИЯВЛЕННЯ АНОМАЛІЙ У ФІНАНСОВИХ ТРАНЗАКЦІЯХ

Зіненко Михайло Євгенович, Харківське територіальне відділення МАН України, Комунальний заклад «Харківський академічний ліцей № 45 Харківської міської ради», 11 клас, м. Харків

Шаповалова Марія Ігорівна, старший викладач кафедри математичного моделювання та інтелектуальних обчислень в інженерії навчально-наукового інституту комп'ютерного моделювання, прикладної фізики та математики Національного технічного університету «Харківський політехнічний інститут», кандидат технічних наук

Арзубов Микола Олексійович, вчитель інформатики Комунального закладу «Харківський академічний ліцей № 45 Харківської міської ради».

З розвитком технологій і сталим економічним зростанням, характерним для сучасного суспільства, шахрайські дії стали більш поширеними у фінансовій сфері, і це щороку обходиться організаціям і споживачам у сотні мільярдів доларів. Шахраї постійно вдосконалюють свої методи, щоб використовувати вразливі місця існуючих заходів протидії, і найбільше потерпає від них фінансовий сектор.

Метою даної роботи є аналіз методів для виявлення фінансових аномалій, їх порівняння та знаходження оптимальних, що дозволить підвищити ступінь захищеності банківських систем. Для вирішення цієї задачі будуть описані та реалізовані деякі популярні методи виявлення викидів.

Об'єктом дослідження є процес здійснення грошових переказів користувачами.

Предметом дослідження є методи виявлення аномальних транзакцій у фінансових системах на основі машинного навчання.

У результаті дослідження буде знайдено один з оптимальних методів виявлення подібних шахрайств, що може сприяти підвищенню ефективності моніторингу та захисту фінансових систем. Ця робота може знайти застосування в подальших дослідженнях у цій сфері, загалом підвищуючи рівень безпеки та мінімізуючи фінансові ризики.

Ключові слова: фінансові транзакції, аналіз даних, логістична регресія, випадковий ліс, градієнтний бустинг, ізоляційний ліс, MLP, машинне навчання

ЗМІСТ

ВСТУП.....	5
РОЗДІЛ 1. ОПИС ПРОБЛЕМИ ФІНАНСОВИХ АНОМАЛІЙ	7
1.1. Загальне визначення аномалій	7
1.2. Види аномалій	8
1.2.1. Точкові аномалії.....	8
1.2.2. Контекстуальні аномалії	8
1.2.3. Колективні аномалії	9
1.3. Проблеми при виявленні відхилень	9
1.3.1. Визначення нормальної області	10
1.3.2. Неоднорідність аномалій	10
1.3.3. Відсутність маркованих даних	10
1.3.4. Наявність шуму в нормальних даних.....	10
1.3.5. Зв'язки між даними.....	10
1.4.Фінансове шахрайство	11
РОЗДІЛ 2. ВИЯВЛЕННЯ АНОМАЛІЙ НА ПРЕДМЕТ ШАХРАЙСТВА.....	12
2.1. Керовані методи виявлення аномалій.....	12
2.2. Некеровані методи виявлення аномалій	12
2.3. Напівкеровані методи виявлення аномалій	13
2.4. Методи виявлення аномалій на основі графів.....	13
2.5. Опис обраних методів для вивчення й порівняння	13
2.5.1. Логістична регресія	13
2.5.2. Випадковий ліс	15
2.5.3. Градієнтний бустинг	16
2.5.4. Ізоляційний ліс	17
2.5.5. Багатошаровий перцептрон Румельхарта	18
РОЗДІЛ 3. АНАЛІЗ ДАНИХ ТА РЕАЛІЗАЦІЯ МЕТОДІВ ВИЯВЛЕННЯ ФІНАНСОВИХ АНОМАЛІЙ	20
3.1. Аналіз набору даних	20
3.1.1. Перевірка відсутніх даних	21
3.1.2. Дисбаланс даних.....	21
3.2. Реалізація алгоритмів та аналіз їх результатів	22
3.2.1. Підготовка набору даних до навчання моделей	22

3.2.2. Навчання моделей	23
3.2.3. Обробка отриманих результатів	24
ВИСНОВКИ	27
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ	29

ВСТУП

Виявлення аномалій, або виявлення викидів, містить у собі виявлення даних або подій, що є відхиленням від очікуваних закономірностей. Для опису таких аномалій використовують різні терміни, такі як неузгоджені об'єкти, винятки, аберації, особливості або забруднення. Хоча методи, які використовуються для виявлення аномалій, можуть бути схожими, вони часто називаються по-різному залежно від конкретної сфери застосування.

Виявлення шаблонів аномальної поведінки або подій має величезне значення, оскільки вони можуть дати суттєві, дієві, а часто й критично важливі відомості в різних галузях. Аномалії можуть виникати з різних причин, включно зі шкідливими діями або збоями в роботі системи, і всі вони становлять інтерес для аналітиків. Наприклад, у даних про кредитні картки аномальна поведінка може свідчити про крадіжку особистих коштів або шахрайські операції, здійснені неавторизованою стороною. Також у структурі мережевого трафіку аномалії можуть свідчити про спробу зловмисників скомпрометувати систему, що може призвести до серйозних збоїв у роботі або вказати на зламаній комп'ютер, який надсилає конфіденційні дані в несанкціоноване місце. У галузі охорони здоров'я методи виявлення аномалій сприяють виявленню злоякісних клітин або ділянок на знімках магнітно-резонансної томографії (МРТ). Крім того, аномальні вимірювання з датчиків космічних апаратів можуть свідчити про несправний компонент, а в природі землетруси можна передбачати, виявляючи аномалії в даних, що їм передують. У всіх цих сферах існує поняття нормальної моделі даних, від якої відхиляються аномалії.

Останнім часом актуальною проблемою стало фінансове шахрайство, що охоплює такі види діяльності, як шахрайство з кредитними картками, страхове шахрайство, відмивання грошей, шахрайство у сфері охорони здоров'я, а також шахрайство з цінними паперами й товарами. Зростаючі масштаби фінансового шахрайства створюють серйозні економічні проблеми, а глобальні втрати

сягають мільярдів доларів на рік. Ці злочини мають далекосяжні наслідки для різних галузей економіки і навіть пов'язані з фінансуванням таких видів незаконної діяльності, як організована злочинність і наркоторгівля. Компанії та продавці несуть фінансовий тягар цих злочинів, зазнаючи витрат, пов'язаних із поверненням платежів, адміністративними витратами та втратою довіри споживачів. З огляду на тяжкі наслідки такого шахрайства, розробка ефективних стратегій і методів його виявлення є нагальною потребою.

Метою роботи є аналіз методів для виявлення фінансових аномалій, їх порівняння та знаходження оптимальних, що дозволить підвищити ступінь захищеності банківських систем. Для вирішення мети будуть описані та реалізовані деякі популярні методи виявлення викидів.

Завданнями роботи є:

– провести детальний аналіз методів виявлення фінансових аномалій на прикладі шахрайства з використанням різних алгоритмів машинного навчання та виокремити потенційно ефективні моделі для їх виявлення.

– використати набір даних, що містить інформацію про транзакції європейських власників кредитних карток.

– застосувати різні типи алгоритмів для виявлення аномалій: логістична регресія, випадковий ліс, градієнтний бустинг, MLP та ізоляційний ліс.

– використати різні метрики для оцінки результатів, такі як точність, влучність, повнота та F1-міра.

Об'єктом дослідження є процес здійснення грошових переказів користувачами.

Предметом дослідження є методи виявлення аномальних транзакцій у фінансових системах на основі машинного навчання.

Науковими результатами дослідження є: подальший розвиток отримали методи виявлення аномалій у фінансових транзакціях, що надало можливість підвищити точність ідентифікації шахрайських операцій шляхом використання різноманітних методів машинного навчання.

РОЗДІЛ 1

ОПИС ПРОБЛЕМИ ФІНАНСОВИХ АНОМАЛІЙ

1.1. Загальне визначення аномалій

Численні автори запропонували різні визначення аномалій, і загальноприйнятого визначення ще не з'явилося. Точна характеристика аномалії залежить від припущень щодо структури даних і конкретного застосування. Проте існують визначення, які вважаються широко застосовними в різних контекстах. Серед них найбільш загальновизнаним є визначення С. Гокінга [1], який визначає аномалію так: "Аномалія - це спостереження, яке настільки відхиляється від інших спостережень, що викликає підозру, що воно було породжене іншим механізмом". Це визначення впливає зі статистичної перспективи, де нормальні дані відповідають певному механізму, тоді як аномалії - це випадки або зразки, які відхиляються від цього механізму. Отже, аномалії часто дають цінну інформацію про нетипові характеристики системи, які впливають на утворювальний механізм.

Виявлення аномалій має схожість із завданням видалення шуму, яке передбачає усунення небажаного шуму з даних. Однак ці два завдання відрізняються одне від одного. У реальних сценаріях на дані може впливати значна кількість шуму, який може не представляти інтересу для аналітика, але є проблемою під час аналізу даних. Як правило, інтерес становлять значні відхилення. Негативний вплив шуму на результати аналізу даних підкреслює важливість його видалення, оскільки це допомагає усунути будь-які небажані елементи перед аналізом.

Аномалії можуть виникати з різних джерел, зокрема людські помилки, несправності приладів, шахрайські дії, поведінкові зміни або несправності всередині системи. Те, як система виявлення відхилень реагує на аномалію, залежить від конкретної сфери застосування. Наприклад, якщо аномалія вказує

на друкарську помилку, допущену працівником, який вводить дані, достатньо простого повідомлення йому для виправлення. Аналогічно, аномальні дані з показань приладів можна просто видалити після їх ідентифікації. Щобільше, у критичних середовищах, таких як системи моніторингу вторгнень або виявлення шахрайства, система виявлення аномалій повинна бути здатна швидко й в реальному часі знаходити відхилення від норми з подачею відповідного сигналу тривоги, щоб забезпечити своєчасне втручання.

1.2. Види аномалій

Аномалії можна класифікувати за трьома різними категоріями, і врахування типу аномалії є критичним фактором для будь-якого методу виявлення аномалій [2].

1.2.1. Точкові аномалії

Точкові аномалії являють собою найпростішу форму аномалії і є основним об'єктом багатьох досліджень у методах виявлення аномалій. Вони характеризуються окремими точками даних або подіями, які значно відхиляються від решти набору даних. У ширшому розумінні, це точки, які виходять за межі встановленої норми або розділювальної межі, прикладом яких є точки *O1* і *O2*. Щоб проілюструвати це на реальному прикладі, розглянемо дані про транзакції з кредитними картками для фізичних осіб, які визначаються єдиною ознакою: сумою покупки. У цьому контексті будь-яка транзакція на суму, що значно перевищує звичайний діапазон витрат для цієї особи, буде класифікована як точкова аномалія.³

1.2.2. Контекстуальні аномалії

Контекстуальні аномалії – це дані, які вважаються аномальними лише в певному контексті, а не інакше. Структура набору даних вводить поняття контексту для цих типів аномалій, і формулювання проблеми вимагає його явної

специфікації. Кожен екземпляр даних визначається двома атрибутами: контекстними та поведінковими. Контекстуальні атрибути визначають контекст для цього екземпляра, наприклад, час у даних часових рядів. Поведінкові атрибути позначають позаконтекстні характеристики, такі як сума покупки в наборі даних про кредитні картки. Аномальна поведінка оцінюється за допомогою цих поведінкових атрибутів у певному контексті. Наприклад, при виявленні шахрайства з кредитними картками контекстним атрибутом може бути час покупки. Якщо припустити, що людина зазвичай витрачає сто доларів на тиждень, за винятком різдвяного тижня, коли вона витрачає тисячу доларів, то покупка на тисячу доларів у середній тиждень у травні вважатиметься контекстуальною аномалією. Це тому, що вона відхиляється від типової поведінки в контексті часу, навіть якщо витратити таку ж суму під час різдвяного тижня вважалося б нормальним явищем.

1.2.3. Колективні аномалії

Колективні аномалії стосуються групи пов'язаних між собою випадків, які колективно відхиляються від норми всього набору даних. Окремі події або дані в межах колективної аномалії не обов'язково можуть бути аномаліями самі по собі; однак їхня поява як сукупності вважається аномальною. Важливо зазначити, що точкові аномалії можуть виникати в будь-якому наборі даних, тоді як колективні аномалії можуть виникати лише в тих наборах даних, де є зв'язок між екземплярами даних. Контекстуальні атрибути можуть бути релевантними лише тоді, коли вони присутні в даних. Крім того, точкові або колективні аномалії можуть навіть вважатися контекстними аномаліями, якщо вони спостерігаються в певному контексті.

1.2. Проблеми при виявленні відхилень

Аномалії, як було визначено раніше, належать до спостережень або подій, які відхиляються від очікуваної нормальної поведінки. Хоча може здатися, що

визначити область, яка представляє нормальну поведінку, і позначити спостереження за її межами як аномалії, може здатися простим, але виявлення аномалій створює надзвичайні складнощі, відмінні від багатьох інших аналітичних і навчальних проблем [3].

1.3.1 Визначення нормальної області

Визначення точної межі, яка б охоплювала всі можливі форми нормальної поведінки, є дуже складним завданням. Межі між нормальною та аномальною поведінкою часто бракує точності. Це може призвести до неправильної класифікації спостережень поблизу цієї межі. Крім того, аномалії, спричинені зловмисною діяльністю, є динамічними та адаптивними, що ускладнює їх виявлення, оскільки зловмисники постійно розвивають свою тактику, щоб імітувати нормальну діяльність.

1.3.2. Неоднорідність аномалій

Аномалії можуть бути гетерогенними, тобто різні класи аномалій можуть мати дуже різні характеристики, що вказують на аномалію. Це додає складності процесу виявлення.

1.3.3. Відсутність маркованих даних

Аномалії, як правило, трапляються рідко порівняно з нормальними випадками, що призводить до дефіциту мічених даних для навчання й перевірки моделей. Проблеми конфіденційності даних і витрати, пов'язані з ручним маркуванням, також сприяють посиленню цієї проблеми.

1.3.4. Наявність шуму в нормальних даних

Наявність шуму, який може діяти подібно до аномалій, ускладнює встановлення чітких меж або правил прийняття рішень у наборі даних.

1.3.5. Зв'язки між даними

Залежно від характеру проблеми, між екземплярами даних можуть існувати взаємозв'язки, наприклад, просторові, часові або графічні дані. Це необхідно враховувати на етапі проектування, щоб переконатися, що обраний метод є надійним і добре працює для конкретного застосування.

1.3. Фінансове шахрайство

Майже кожен набір даних має нерівномірне представлення класів. Це не є проблемою, якщо різниця невелика. Однак, коли один або декілька класів є дуже нечисленними, більшість моделей погано працюють при визначенні класів, що становлять меншість. Однією з таких областей є фінансове шахрайство [4].

За останні десять років фінансове шахрайство привернуло до себе дуже багато уваги через потенційні наслідки невиявлених аномалій для промисловості та повсякденного життя. Загальноприйнятого визначення не існує, але найближчим є визначення Асоціації сертифікованих експертів з питань шахрайства: «будь-яка навмисна дія, спрямована на позбавлення іншої особи власності або грошей шляхом підступності, обману або інших нечесних засобів».

Зі стрімким поширенням і розвитком сучасних технологій, таких як мережа «Інтернет», дротових пристроїв, телефони й ноутбуки, а також соціальних мереж, зросла й шахрайська активність. Це призвело до мільярдних збитків бізнесу, що, відповідно, спонукало до активних зусиль з вивчення методів виявлення аномалій для викриття шахрайства.

У цьому дослідженні ми розглядали застосування методів виявлення аномалій у шахрайствах з кредитними картками.

РОЗДІЛ 2

ВИЯВЛЕННЯ АНОМАЛІЙ НА ПРЕДМЕТ ШАХРАЙСТВА

Виявлення аномалій сприяє викриттю фінансового шахрайства й використовується для вилучення інформації з великих обсягів даних. Часто методи класифікуються за чотирма ознаками: керовані, некеровані, напівкеровані та методи на основі графів [5].

2.1. Керовані методи виявлення аномалій

При контрольованому виявленні аномалій алгоритм навчається на маркованому наборі даних, де чітко визначені як нормальні, так й аномальні екземпляри. Поширені методи включають методи опорних векторів (SVM), k -найближчих сусідів (KNN), дерева рішень та ансамблеві методи. Під час навчання модель вивчає зразки нормальної поведінки, а потім може ідентифікувати аномалії на основі відхилень від цієї вивченої нормальної поведінки.

2.2. Некеровані методи виявлення аномалій

Неконтрольоване виявлення аномалій виконується на наборах даних, де для навчання доступні лише нормальні екземпляри. Алгоритми кластеризації (наприклад, кластеризація за методом k -середніх), ізоляційні ліси та однокласові SVM є поширеними неконтрольованими методами. Ці алгоритми спрямовані на виявлення випадків, які відхиляються від загальних закономірностей, що спостерігаються в даних, припускаючи, що аномалії є рідкісними й мають чіткі характеристики.

2.3. Напівкеровані методи виявлення аномалій

Напівкеровані методи знаходяться між керованими й некерованими підходами. Вони використовують набір даних, що містить переважно нормальні екземпляри й меншу підмножину позначених аномалій.

Для напівкерованого виявлення аномалій можна використовувати такі методи, як автокодера (тип нейронної мережі).

2.4. Методи виявлення аномалій на основі графів

Графові методи використовують взаємозв'язки між точками даних для виявлення аномалій. Виявлення аномалій на основі графів часто використовується в сценаріях, де дані можуть бути представлені у вигляді графа, таких як соціальні мережі, або виявлення мережевих вторгнень. Для виявлення аномалій на основі графів можна адаптувати такі алгоритми, як ліс ізоляції на графах, спектральна кластеризація та локальний фактор відхилення (LOF).

2.5. Опис обраних методів для вивчення й порівняння

Поширеною проблемою в задачах виявлення шахрайства є те, що дані часто дуже незбалансовані [6]. Це означає, що кількість шахрайських транзакцій значно менша, ніж кількість законних транзакцій. У таких випадках звичайні моделі машинного навчання, як правило, упереджено ставляться до класу більшості й можуть погано виявляти клас меншості.

2.5.1. Логістична регресія

Логістична регресія [7] є лінійним методом, який добре працює для бінарної класифікації. Метод є ефективним в задачах виявлення шахрайства через те, що застосовується у випадку, коли залежна змінна може набувати тільки двох значень (0 або 1).

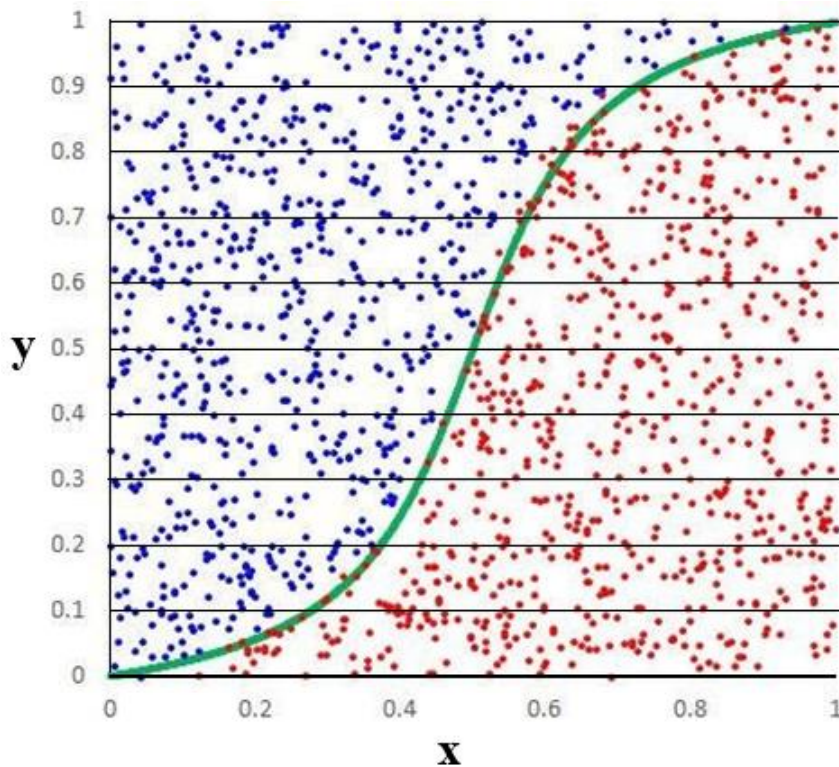


Рис. 2.1. Логістична регресія (рисунок автора)

На початку будується лінійна комбінація вхідних ознак з їхніми вагами.

Після чого лінійна комбінація z піддається логістичній функції, також відомій як сигмоїда, щоб перетворити результат у діапазон від 0 до 1.

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

де e – основа натурального логарифма.

Ця функція бере значення лінійної комбінації та перетворює його на ймовірність. Якщо ймовірність наближена до 1, то модель вважає, що об'єкт належить до позитивного класу, а якщо до 0, то до негативного.

2.5.2. Випадковий ліс

Випадковий ліс [8] - це ансамбль дерев рішень, який може бути ефективним в обробці складних, нелінійних залежностей, який працює за допомогою побудови численних дерев прийняття рішень.

Для побудови кожного дерева в лісі використовується підмножина навчальної вибірки, яка утворюється шляхом випадкового вибору з поверненням. Під час побудови кожного вузла дерева з випадкової підмножини ознак вибирають найкращу ознаку для поділу.

Для кожного дерева використовується стандартний алгоритм побудови дерева рішень такий як *Classification and Regression Trees (CART)*. Дерев будуються доти, доки не досягнуто певної глибини або не виконано критерії зупинки.

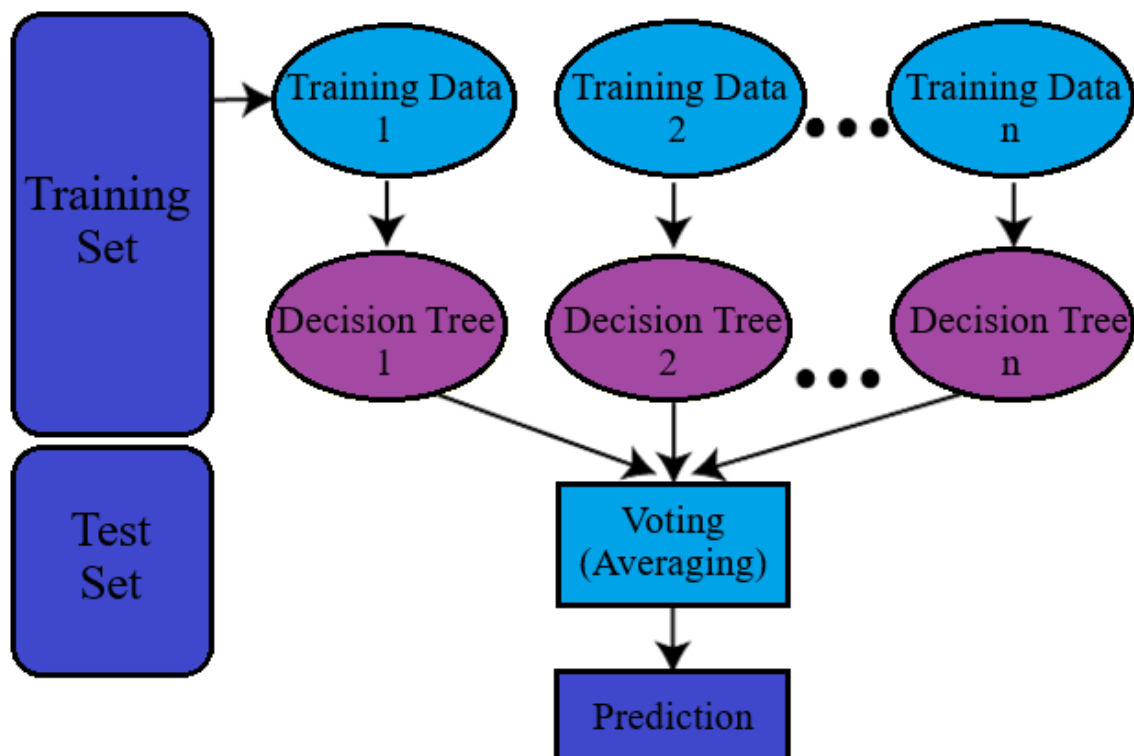


Рис. 2.2. Випадковий ліс (рисунок автора)

Коли всі дерева побудовані, передбачення кожного дерева об'єднуються для отримання остаточного рішення. У разі класифікації ухвалюється рішення більшості дерев, а в разі регресії береться середнє значення передбачень.

2.5.3. Градієнтний бустинг

Градієнтний бустинг [9] - це ансамбль, який комбінує кілька слабких моделей (зазвичай дерев рішень) для поліпшення передбачуваної здатності.

Зазвичай як базову модель використовують дерева рішень невеликої глибини. Вони є слабкими моделями, оскільки можуть давати обмежену якість передбачень.

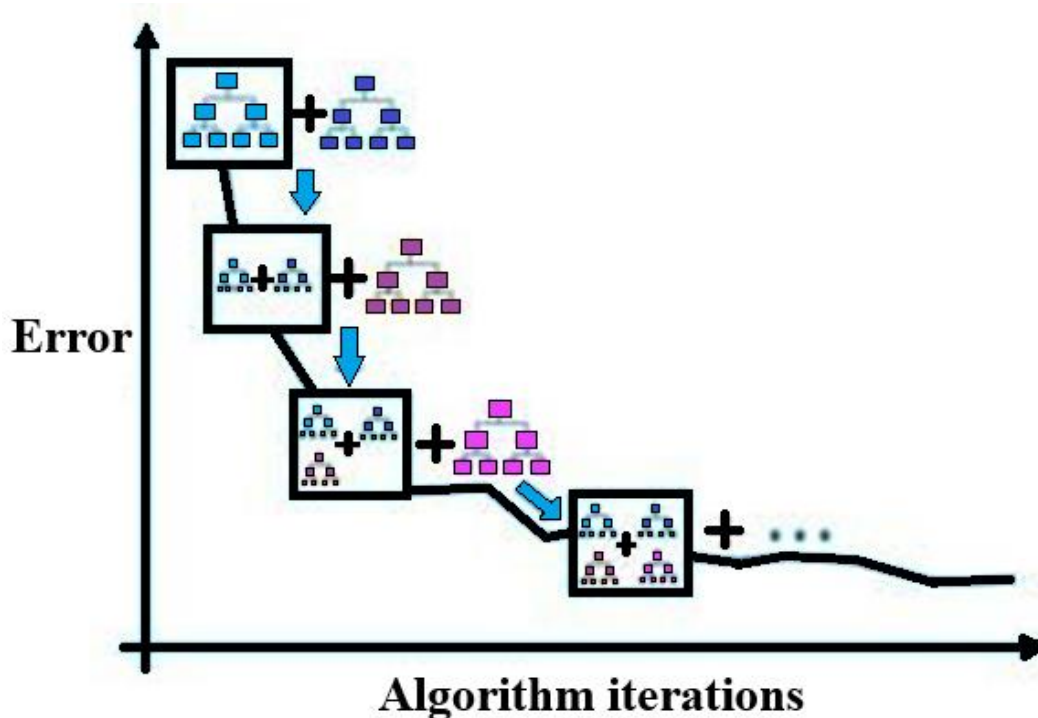


Рис. 2.3. Градієнтний бустинг (рисунок автора)

Моделі навчаються ітеративно. На кожній ітерації будується нова модель, яка зосереджується на помилках попередніх моделей. Це відбувається через підстроювання параметрів моделі так, щоб вона передбачала рештки попередніх моделей.

Важливим елементом градієнтного бустингу є визначення функції втрат, яка вимірює помилку передбачення.

У побудові градієнтного бустингу було прийнято рішення використовувати метод eXtreme Gradient Boosting (XGBoost), який на відміну від стандартного градієнтного бустингу, додає декілька оптимізацій:

- можливість паралельного навчання на рівні дерев, що прискорює процес;
- використання спеціальної структури даних, яка прискорює пошук найкращого поділу в деревах;
- автоматична обробка пропущених значень в даних.

2.5.4. Ізоляційний ліс

Ізоляційний ліс [10] – це метод, заснований на аналізі того, наскільки швидко об'єкт можна виокремити "ізолювати" в дереві.

Будується набір дерев. Для кожного дерева в лісі вибираються випадкові підвибірки даних. Кожне дерево будується шляхом рекурсивного поділу даних. На кожному рівні дерева вибирають випадкову ознаку, і дані поділяють за пороговим значенням цієї ознаки. Цей процес повторюється, поки кожен екземпляр даних не буде ізолюваний у своєму вузлі.

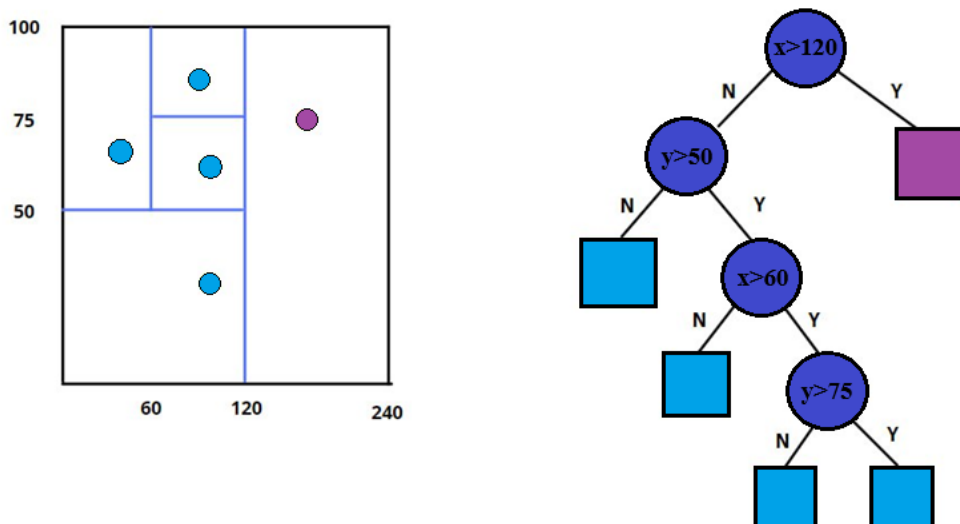


Рис 2.4. Ізоляційний ліс (рисунок автора)

Кількість поділів необхідних для ізоляції кожної точки, використовується як міра «ізоляції». Аномалії, як правило, потребуватимуть менше поділів, ніж нормальні точки.

Рішення про наявність аномалії: За результатами вимірювання ізоляції для кожного дерева вирішується, чи є точка аномалією. Це може бути зроблено, наприклад, шляхом визначення порога ізоляції, нижче за який точка вважається аномалією

2.5.5. Багатошаровий перцептрон Румельхарта

Багатошаровий перцептрон Румельхарта (MLP) [11] - це різновид штучної нейронної мережі, що складається з декількох шарів нейронів, включно з вхідним шаром, одним або кількома прихованими шарами та вихідним шаром.

На вході MLP знаходиться вектор ознак, що представляє вхідні дані. Кожен елемент цього вектора являє собою одну ознаку.

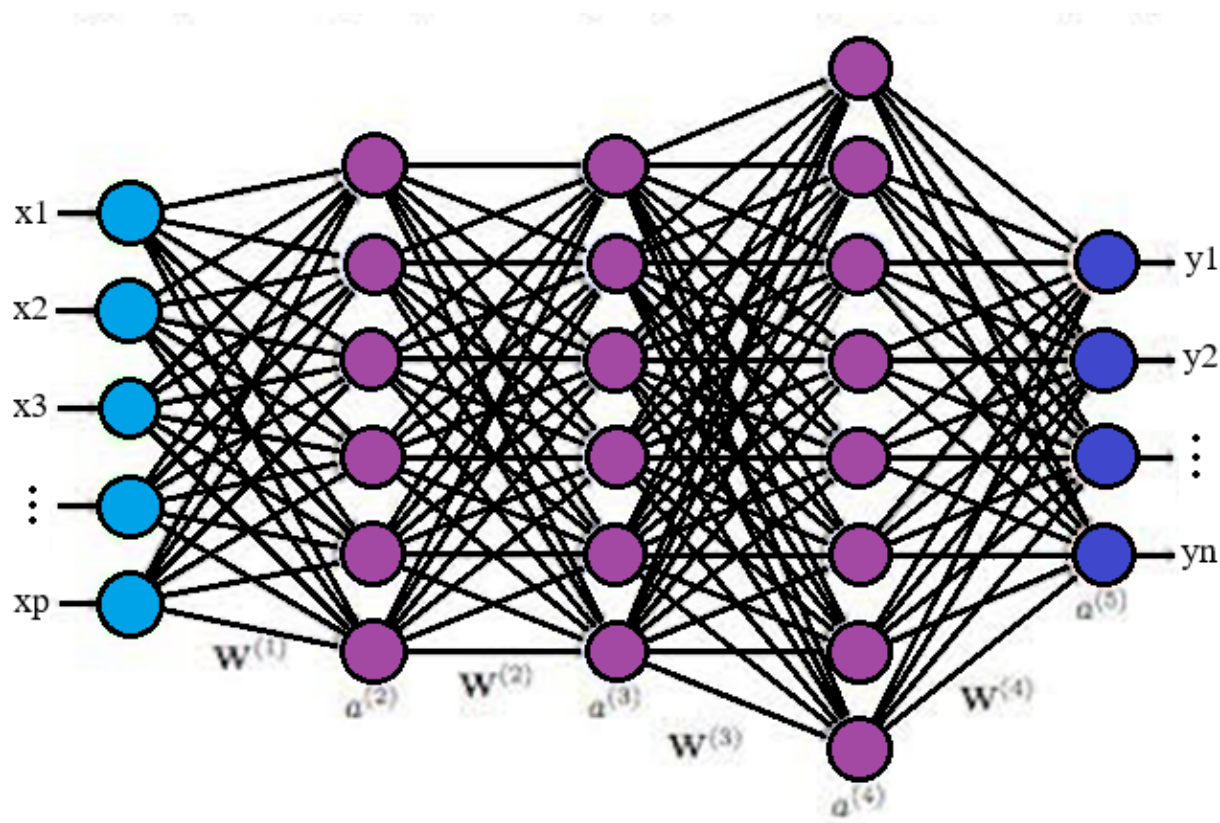


Рис 2.5. Багатошаровий перцептрон Румельхарта (рисунок автора)

Після вхідного шару йдуть один або кілька прихованих шарів. Кожен нейрон у прихованому шарі з'єднаний з усіма нейронами попереднього шару. Кожен зв'язок між нейронами має свою вагу, яка визначає важливість цього зв'язку. Нейрони в прихованих шарах приймають зважені суми вхідних сигналів

і застосовують активаційну функцію до цієї суми. Активаційна функція допомагає введенню нелінійності в мережу, що дає змогу MLP моделювати складніші залежності.

Останній прихований шар з'єднаний із вихідним шаром. Нейрони у вихідному шарі аналогічно приймають зважені суми сигналів від попереднього шару й застосовують активаційну функцію.

MLP навчається з використанням методів зворотного поширення помилки. Цей процес передбачає обчислення помилки між передбаченим значенням і фактичним значенням, а потім коригування ваг мережі таким чином, щоб зменшити цю помилку.

РОЗДІЛ 3

АНАЛІЗ ДАНИХ ТА РЕАЛІЗАЦІЯ МЕТОДІВ ВИЯВЛЕННЯ ФІНАНСОВИХ АНОМАЛІЙ

Робота з набором даних та реалізація методів, використаних у дослідженні, проводилася у середовищі Google Colab (<https://bit.ly/48I51Co>). Звідти також були отримані діаграми, таблиці та блоки коду.

3.1. Аналіз набору даних

Набір даних містить транзакції, здійснені європейськими власниками кредитних карток у вересні 2013 року. У цій вибірці представлені транзакції, що відбулися протягом двох днів, де ми маємо 492 шахрайства з 284 807 транзакцій. Набір даних дуже незбалансований, позитивний клас (шахрайство) становить 0,173% від усіх транзакцій.

Навчальна вибірка містить лише числові вхідні змінні, які є результатом PCA-перетворення. Ознаки V1, V2, ... V28 - це головні компоненти, отримані за допомогою PCA. Єдині ознаки, які не були перетворені за допомогою PCA, - це Time та Amount. Характеристика Time містить секунди, що пройшли між кожною транзакцією та першою транзакцією в наборі даних. Ознака Amount - це сума транзакції. Клас ознаки - це змінна відгуку, яка приймає значення 1 у випадку шахрайства та 0 у протилежному випадку.

	Time	V1	V2	V3	V4	V5	V6	V7	\
0	0.0	-1.359807	-0.072781	2.536347	1.378155	-0.338321	0.462388	0.239599	
1	0.0	1.191857	0.266151	0.166480	0.448154	0.060018	-0.082361	-0.078803	
2	1.0	-1.358354	-1.340163	1.773209	0.379780	-0.503198	1.800499	0.791461	
3	1.0	-0.966272	-0.185226	1.792993	-0.863291	-0.010309	1.247203	0.237609	
4	2.0	-1.158233	0.877737	1.548718	0.403034	-0.407193	0.095921	0.592941	

	V8	V9	...	V21	V22	V23	V24	V25	\
0	0.098698	0.363787	...	-0.018307	0.277838	-0.110474	0.066928	0.128539	
1	0.085102	-0.255425	...	-0.225775	-0.638672	0.101288	-0.339846	0.167170	
2	0.247676	-1.514654	...	0.247998	0.771679	0.909412	-0.689281	-0.327642	
3	0.377436	-1.387024	...	-0.108300	0.005274	-0.190321	-1.175575	0.647376	
4	-0.270533	0.817739	...	-0.009431	0.798278	-0.137458	0.141267	-0.206010	

	V26	V27	V28	Amount	Class
0	-0.189115	0.133558	-0.021053	149.62	0
1	0.125895	-0.008983	0.014724	2.69	0
2	-0.139097	-0.055353	-0.059752	378.66	0
3	-0.221929	0.062723	0.061458	123.50	0
4	0.502292	0.219422	0.215153	69.99	0

Рис. 3.1. Перші рядки вибірки даних (таблиця автора)

3.1.1. Перевірка відсутніх даних

Першим етапом попередньої обробки даних є перевірка на порожні комірки в навчальній вибірці. Для цього, необхідно розробити скрипт, який дозволяє перевірити кожну ознаку кожного об'єкту на відсутність інформації в них. На рис. 3.2 представлено скрипт, який дозволяє виконати цей етап обробки даних.

```
pd.DataFrame(data.isnull().sum(), columns=["Missing Values"])
```

Рис. 3.2. Перевірка на порожні комірки (авторський блок коду)

Цей скрипт створює DataFrame такого ж розміру, що й початкові дані, але з булевими значеннями True для комірок із пропущеними значеннями і False або NaN для заповнених комірок. Після цього підраховується сума True значень, що і дає кількість пропущених значень у кожному стовпчику.

3.1.2. Дисбаланс даних

Перевіримо незбалансованість даних відносно до цільового значення, тобто класу. На рис. 3.3 лише 492 (або 0,173%) транзакції є шахрайськими. Це означає, що дані дуже незбалансовані щодо цільової змінної.



Рис. 3.3.. Розподіл цільової змінної (діаграма автора)

Ми не будемо використовувати техніки *undersampling* чи *oversampling*, бо це призводить до втрати важливої інформації та зниженню продуктивності моделей. Вони можуть переоцінювати кількість шахрайських транзакцій і, зрештою, помилково визначати багато законних транзакцій.

3.2. Реалізація алгоритмів та аналіз їх результатів

3.2.1. Підготовка набору даних до навчання моделей

Для навчання наших моделей потрібно зробити перетворення з початковим набором даних. Для цього поділимо дані на ознаки (X) і цільову змінну (Y). Після цього розіб'ємо вибірку даних на навчальний і тестовий набори. Використаємо для цього скрипт представлений на рис. 3.4.

```
X = data.drop('Class', axis=1)
y = data['Class']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Рис. 3.4. Розбиття вибірки на навчальний та тестовий набори
(авторський блок коду)

Перед навчанням моделей також часто потрібне масштабування даних. Це важливо, оскільки багато алгоритмів чутливі до масштабу ознак. Скористуємося стандартизацією, це перетворення ознак так, щоб вони мали середнє значення 0 і стандартне відхилення 1 Використаємо *StandardScaler* з бібліотеки *Python scikit-learn*, реалізація відображена на рис. 3.5

```
scaler = StandardScaler()  
X_scaled = scaler.fit_transform(X_train)  
X_test_scaled = scaler.transform(X_test)
```

Рис. 3.5. Масштабування даних (авторський блок коду)

3.2.2. Навчання моделей

Після підготовки вибірки даних, ініціалізуємо кожену модель задаючи гіперпараметри, для всіх моделей це *random_state*, цей параметр використовується для встановлення початкового стану генератора випадкових чисел.

Для ізоляційного лісу ми встановлюємо параметр *contamination*, який визначає очікувану долю аномалій в даних. У нашому випадку це 0,173.

А для нейронної мережі – параметр *max_iter*, який визначає максимальну кількість епох, які модель MLP виконуватиме під час навчання. Навчання може завершитися раніше, якщо сходиться до оптимуму.

Проведемо навчання кожної з обраних моделей. Узагальнений код представлено на рис. 3.6

```
model.fit(X_scaled, y_train)  
y_pred = model.predict(X_test_scaled)
```

Рис. 3.6. Узагальнений код для тренування моделей (авторський блок коду)

Змінна *model* – об'єкт моделі. Перший рядок відповідає за навчання, *X_scaled* - це матриця ознак, яка була попередньо масштабована, а *y_train* – це

вектор цільових значень для навчання моделі. У процесі виконання модель налаштовує свої внутрішні параметри, щоб краще відповідати набору даних.

Другий рядок відповідає за передбачення даних, після того як модель вже навчена. *X_text_scaled* – це матриця масштабованих ознак тестового набору даних, а *y_pred* – вектор передбачених значень, отриманих моделлю на основі тестових даних.

Тільки для ізоляційного лісу треба перевести передбачення в бінарні значення. Тому що передбачення ізоляційного лісу представлені значенням -1 для аномальних і 1 для нормальних даних. Переведення в бінарні значення спростить відображення результатів. Скрипт для виконання цієї операції приведений на рис. 3.7

```
y_pred_iso_forest_binary = [1 if x == -1 else 0 for x in y_pred_iso_forest]
```

Рис. 3.7 – Переведення передбачень у бінарні значення (авторський блок коду)

3.2.3. Обробка отриманих результатів

Через те що вибірка даних незбалансована, тобто 99,8% операцій не є шахрайськими, модель, яка правильно прогнозує, що всі транзакції не є шахрайськими, досягне точності 99,8%. Зважаючи на це, можемо зробити висновок, що точність сама по собі не є достатнім способом оцінки роботи моделей. Тому представимо результати як набір кількох метрик. Спочатку подивимось на абстрактну матрицю невідповідностей представлену в рис 3.8

		Прогнозований клас	
		Predicted Positive (PP)	Predicted Negative (PN)
Справжній клас	Загальна кількість = P + N		
	Positive (P)	True positive (TP)	False negative (FN)
	Negative (N)	False positive (FP)	True negative (TN)

Рис. 3.8. Матриця невідповідностей (рисунок взятий з <http://bit.ly/49D4Xq1>)

У матриці ми бачимо чотири компоненти:

- істинно позитивні (TP) – випадки, коли модель правильно ідентифікувала шахрайські транзакції,
- хибно позитивні (FP) – випадки, коли модель помилково визначила, що нешахрайські транзакції є шахрайськими,
- істинно негативні (TN) – випадки, коли модель правильно ідентифікувала нешахрайські транзакції,
- хибно негативні (FN) – випадки, коли модель помилково визначила, що шахрайські транзакції не є шахрайськими.

Тепер розглянемо кожну з метрик [12], за якими буде проводитися оцінка моделей. Оцінювати модель будемо за результатами виявлення шахрайських операцій.

1. Точність – вимірює частку вірних передбачень відносно до всіх передбачень. Відношення кількості правильно класифікованих прикладів (як шахрайські, так і нешахрайські транзакції) до загальної кількості прикладів.
2. Влучність – є часткою релевантних зразків серед знайдених. Відношення кількості правильно класифікованих шахрайських транзакцій до загальної кількості транзакцій, які були класифіковані як шахрайські.
3. Повнота – є часткою загального числа позитивних зразків, яку було дійсно знайдено. Відношення кількості правильно класифікованих шахрайських транзакцій до загальної кількості фактично шахрайських транзакцій.
4. F1-міра – це середньозважена гармонійна між влучністю й повнотою. Використовується, коли важливі і влучність, і повнота.

Отримані результати для кожної моделі наведені в табл. 3.1. Кольором позначено наскільки добре модель впоралася за даною метрикою, зелений колір – найкраще, червоний колір – найгірше.

Таблиця 3.1

Результати навчання моделей

	Точність	Влучність	Повнота	F1-міра
Логістична регресія	99,91%	86%	58%	70%
Випадковий ліс	99,96%	97%	77%	86%
Гرادієнтний бустинг	99,96%	96%	78%	86%
Ізоляційний ліс	82,98%	1%	93%	2%
MLP	99,96%	88%	83%	85%

Вибір моделі для задач класифікації залежить від сформованих вимог. Якщо пріоритетом є максимальне виявлення шахрайства, але при цьому допускається певна кількість хибних спрацьовувань, ми можемо обрати модель з високим значенням повноти. Якщо ж пріоритетом є уникнення помилкових звинувачень за рахунок пропуску деяких шахрайств, то краще обрати високоточну модель. При збалансованому підході обирається модель з оптимальними співвідношенням між повнотою та влучністю. Аналіз результатів показав, що ефективними методами виявлення шахрайських транзакцій є випадковий ліс та градієнтний бустинг. Найвища влучність, що означає найнижче невдоволення клієнтів через скасування їх транзакцій, та досить висока повнота, тобто кількість пропущених шахрайських транзакцій – низька, відносно загальної кількості шахрайств.

Однак варто відзначити метод ізоляційного лісу, який, незважаючи на занадто низьку влучність та F1-міру, має найвищий показник повноти, що свідчить про те, що більшість шахрайських транзакцій були виявлені.

ВИСНОВКИ

1. У роботі проведено детальний аналіз методів виявлення фінансових аномалій на прикладі шахрайства з використанням різних алгоритмів машинного навчання.

2. Використано набір даних, що містить інформацію про транзакції європейських власників кредитних карток. Набір даних був дуже незбалансований, з лише 0,173% шахрайських транзакцій. Використано числові ознаки, які були результатом PCA-перетворення, та додатково включено час і суму транзакції.

3. Застосовано різні типи алгоритмів для виявлення аномалій: логістична регресія, випадковий ліс, градієнтний бустинг, MLP та ізоляційний ліс.

4. Використано різні метрики для оцінки результатів, такі як точність, влучність, повнота та F1-міра.

5. Кожен алгоритм має свої переваги та обмеження в залежності від конкретного випадку використання.

- Логістична регресія є простою й легко інтерпретованою, слугує базовою моделлю для порівняння зі складнішими методами. Показала найнижчі, але непогані результати.
- Випадковий ліс ефективний метод в обробці складних, нелінійних залежностей. За результатами дослідження є один з найефективніших методів виявлення шахрайських транзакцій. Має високі показники влучності та F1-міри.
- Градієнтний бустинг має результати схожі з випадковим лісом. Висока влучність і повнота свідчать про ефективність моделі для даної навчальної вибірки.
- Ізоляційний ліс допускає багато помилок першого типу. Але найвищий показник повноти – 0,93, що свідчить про те, що порівняно з іншими

методами, ізоляційний ліс найкраще впорався з виявленням шахрайських транзакцій.

- Багатошаровий перцептрон Румельхарта має усереднені показники, дещо знижена влучність, але високий показник повноти, дозволяє досягнути F1-міри, порівняної з випадковим лісом і градієнтним бустингом.

6. Проведений аналіз дозволяє виокремити потенційно ефективні моделі для виявлення фінансових аномалій. Отримані результати дозволяють зробити висновок про можливість використання методів класифікації випадкового лісу, градієнтного бустингу, багатошарового перцептрон Румельхарта та ізоляційного лісу в якості основних та додаткових засобів виявлення аномалій у фінансових транзакціях. Подальшим розвитком наукової роботи може бути об'єднання розглянутих алгоритмів в окремий ансамблевий метод типу стекінг.

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Hawkins, Douglas M. Identification of outliers. Vol. 11. London: Chapman and Hall, 1980.
2. Auxeno. “ Fraud Detection: Tackling Imbalanced Data.” Kaggle, 13 May 2023, URL: www.kaggle.com/code/auxeno/fraud-detection-tackling-imbalanced-data. (дата звернення: 06.03.2024)
3. Gpreda. “Credit Card Fraud Detection Predictive Models.” Kaggle, 1 Apr. 2021, URL: www.kaggle.com/code/gpreda/credit-card-fraud-detection-predictive-models/notebook. (дата звернення: 06.03.2024)
4. Janiobachmann. “Credit Fraud || Dealing with Imbalanced Datasets.” Kaggle, 3 July 2019, URL: www.kaggle.com/code/janiobachmann/credit-fraud-dealing-with-imbalanced-datasets/notebook. (дата звернення: 06.03.2024)
5. Anomaly detection in Financial Data - Proceedings of the 2023 2nd International Conference on Economics, Smart Finance and Contemporary Trade (ESFCT 2023) (pp.419-426) – Yetong Li
6. Bahnsen, Alejandro Correa, et al. "Feature engineering strategies for credit card fraud detection." Expert Systems with Applications 51 (2016): 134-142.
7. Pant, Ayush. “Introduction to Logistic Regression.” Medium, Towards Data Science, 22 Jan. 2019, URL: towardsdatascience.com/introduction-to-logistic-regression-66248243c148. (дата звернення: 06.03.2024)
8. Deng, Houtao. “An Introduction to Random Forest.” Medium, Towards Data Science, 26 Apr. 2021, URL: towardsdatascience.com/random-forest-3a55c3aca46d. (дата звернення: 06.03.2024)
9. Brownlee, Jason. “A Gentle Introduction to the Gradient Boosting Algorithm for Machine Learning.” MachineLearningMastery.Com, 14 Aug. 2020, URL: machinelearningmastery.com/gentle-introduction-gradient-boosting-algorithm-machine-learning/. (дата звернення: 06.03.2024)

10. Mavuduru, Amol. “How to Perform Anomaly Detection with the Isolation Forest Algorithm.” Medium, Towards Data Science, 8 Apr. 2022, towardsdatascience.com/how-to-perform-anomaly-detection-with-the-isolation-forest-algorithm-e8c8372520bc. (дата звернення: 06.03.2024)
11. Bento, Carolina. “Multilayer Perceptron Explained with a Real-Life Example and Python Code: Sentiment Analysis.” Medium, Towards Data Science, 30 Sept. 2021, URL: towardsdatascience.com/multilayer-perceptron-explained-with-a-real-life-example-and-python-code-sentiment-analysis-cb408ee93141. (дата звернення: 06.03.2024)
12. “Sklearn.Metrics.Classification_report.” Scikit, URL: scikit-learn.org/stable/modules/generated/sklearn.metrics.classification_report.html. (дата звернення: 06.03.2024)