

ОДЕСЬКИЙ НАЦІОНАЛЬНИЙ УНІВЕРСИТЕТ ІМЕНІ І. І. МЕЧНИКОВА

(повне найменування закладу вищої освіти)

Факультет математики, фізики та інформаційних технологій

(повне найменування факультету)

Кафедра інформаційних технологій

(повна назва кафедри)

Кваліфікаційна робота

на здобуття ступеня вищої освіти «Бакалавр»

**«Розробка застосунку для кластеризації даних і
виявлення аномалій»**

(тема кваліфікаційної роботи українською мовою)

**«Development of an application for data clustering and
anomaly detection»**

(тема кваліфікаційної роботи англійською мовою)

Виконала: здобувачка денної форми навчання
спеціальності 122 Комп'ютерні науки
(код, назва спеціальності)

Освітня програма Комп'ютерні науки
(назва)

ПОЛЯК Яна Альфівна

(прізвище, ім'я, по-батькові здобувача)

Керівник к.т.н., доцент Фразе-Фразенко О.О. _____
(науковий ступінь, вчене звання, прізвище, ініціали) (підпис)

Рецензент к.т.н., начальник ІОЦ ОНЕУ, Домаскін О.М.
(науковий ступінь, вчене звання, прізвище, ініціали)

Рекомендовано до захисту:
Протокол засідання кафедри
Інформаційних технологій

№ _____ від _____ 2025 р.

Завідувачка кафедри

КАЗАКОВА Надія

(підпис)

(прізвище, ім'я)

Захищено на засіданні ЕК № _____
протокол № ____ від _____ 2025 р.

Оцінка _____ / _____ / _____
(за національною шкалою/шкалою ECTS/ бали)

Голова ЕК

КОПИЧЕНКО Іван

(підпис)

(прізвище, ім'я)

Одеса 2025

АНОТАЦІЯ

У бакалаврській кваліфікаційній роботі здійснено розробку застосунку для кластеризації даних і виявлення аномалій, орієнтованого на роботу з багатовимірними наборами даних у різних галузях.

Метою роботи є створення консольного програмного продукту на мові Python, що дозволяє автоматизувати процеси попередньої обробки даних, виявлення структурних закономірностей методом кластеризації та визначення аномальних спостережень за допомогою сучасних методів аналізу.

У межах роботи проведено аналіз предметної області, обґрунтовано вибір мов програмування, бібліотек та середовища розробки. Основну увагу приділено алгоритмам K-Means, DBSCAN, Hierarchical Clustering та методу Local Outlier Factor.

Розроблена архітектура застосунку є модульною, включає засоби обробки даних, аналітичні обчислення, візуалізацію результатів і підтримку сценаріїв збереження результатів у форматах CSV та PDF.

Проведено тестування на синтетичних та реальних наборах: дані трафіку, фінансові показники, IoT-дані. Отримано графіки, що підтверджують ефективність роботи застосунку. Результати кластеризації та аномалій наочно демонструють можливості запропонованого рішення.

ABSTRACT

This bachelor's qualification work presents the development of an application for data clustering and anomaly detection, aimed at working with multidimensional datasets across various domains.

The purpose of the project is to create a console-based software product in Python that automates data preprocessing, uncovers structural patterns using clustering algorithms, and identifies anomalous observations with modern analytical techniques.

The research includes analysis of the subject area, justification of the chosen programming language, libraries, and development environment. Particular focus is given to the K-Means, DBSCAN, and Hierarchical Clustering algorithms, as well as the Local Outlier Factor method.

The designed application architecture is modular and incorporates data handling, analytical computations, result visualization, and the ability to export outcomes in CSV and PDF formats.

The system was tested on both synthetic and real-world datasets, including traffic logs, financial indicators, and IoT telemetry. The resulting charts and outputs demonstrate the application's efficiency and the accuracy of clustering and anomaly detection results.

ЗМІСТ

ЗМІСТ	5
ВСТУП	6
1 АНАЛІЗ ЗАДАЧ КЛАСТЕРИЗАЦІЇ ТА ВИЯВЛЕННЯ АНОМАЛІЙ	7
1.1. Постановка задачі аналізу даних	7
1.2. Кластеризація	8
1.3. Виявлення аномалій.....	14
1.4. Огляд існуючих рішень	17
2 ОБҐРУНТУВАННЯ ВИБОРУ ТЕХНОЛОГІЙ І ІНСТРУМЕНТІВ	22
2.1. Огляд мов програмування для аналізу даних	22
2.2. Бібліотеки та програмні засоби для реалізації застосунку	29
3 ПРОЄКТУВАННЯ ЗАСТОСУНКУ	36
3.1. Архітектура застосунку	36
3.2. Вимоги до застосунку.....	38
3.3. Логіка роботи та структура модулів	41
4 РЕАЛІЗАЦІЯ ТА ТЕСТУВАННЯ ЗАСТОСУНКУ	45
4.1. Загальні принципи реалізації	45
4.2. Опис реалізації основних компонентів.....	49
4.3. Тестування та результати	52
ВИСНОВКИ.....	57
ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАННЯ	59

ВСТУП

У сучасному світі обсяг даних, які генеруються в результаті діяльності компаній, установ, соціальних мереж, пристроїв Інтернету речей (IoT) та інших джерел, зростає експоненціально. Це вимагає ефективних методів обробки, аналізу та візуалізації даних. Зокрема, дедалі більшої актуальності набувають задачі автоматичної класифікації, сегментації та виявлення нетипових або підозрілих записів у великих масивах інформації. Такі завдання виникають у найрізноманітніших сферах: від кібербезпеки до фінансового моніторингу, від охорони здоров'я до управління виробничими процесами.

Кластеризація як метод аналізу даних дозволяє розділити об'єкти на групи за спільними ознаками без попереднього маркування. Вона є основою для вирішення багатьох практичних задач, таких як сегментація клієнтів, виділення тематичних груп у текстових даних, агрегація сенсорної інформації тощо. Водночас виявлення аномалій (аномалійний аналіз) дозволяє виявляти об'єкти або події, що суттєво відрізняються від більшості даних, і тому можуть бути індикаторами помилок, збоїв, шахрайства чи нових закономірностей.

У зв'язку з цим виникає потреба у створенні прикладного програмного засобу, який дозволить реалізовувати основні алгоритми кластеризації та виявлення аномалій для аналізу реальних даних.

Об'єктом дослідження є процеси обробки, групування та аналізу даних.

Предметом дослідження є алгоритми кластеризації та методи виявлення аномалій, а також способи їх програмної реалізації.

У процесі дослідження використовувалися методи машинного навчання, теорії кластерного аналізу, чисельні методи, програмна реалізація на мові Python, а також методи моделювання та тестування.

Робота містить 60 сторінок, 19 рисунків, 5 таблиць, 14 джерел посилань

1 АНАЛІЗ ЗАДАЧ КЛАСТЕРИЗАЦІЇ ТА ВИЯВЛЕННЯ АНОМАЛІЙ

1.1. Постановка задачі аналізу даних

Зі стрімким зростанням обсягів цифрової інформації, яка генерується в різних галузях — від телекомунікацій і промислових систем до медицини, фінансів і соціальних мереж — аналіз даних перетворився на ключовий інструмент підтримки прийняття рішень [3]. Сучасні організації активно впроваджують підходи, що дозволяють отримувати нові знання зі складних, багатовимірних наборів даних, часто без чіткого попереднього уявлення про внутрішню структуру цих даних. У цьому контексті особливо важливими є задачі кластеризації — автоматичного групування об'єктів за схожими ознаками — та виявлення аномалій, тобто пошуку нетипових або виняткових спостережень, що можуть містити цінну інформацію [6], [10].

Задача кластеризації полягає в тому, щоб, маючи множину об'єктів без попереднього маркування (тобто без навчальних міток), розділити їх на кластери так, щоб об'єкти в межах одного кластеру були максимально подібними між собою, а об'єкти з різних кластерів — істотно відрізнялися [3], [7]. Це дозволяє автоматично виявляти приховані структури в даних, що особливо корисно на початкових етапах аналізу або коли відсутні дані для навчання. На відміну від класифікації, яка є методом контрольованого навчання і потребує навчального набору з відповідями, кластеризація є методом неконтрольованого навчання, що робить її придатною в ситуаціях із мінімальною або відсутньою апріорною інформацією про структуру даних [6].

Задача виявлення аномалій, яка також відноситься до неконтрольованого або напівконтрольованого навчання, полягає у виявленні об'єктів або записів, що значно відхиляються від загального розподілу [5]. Такі об'єкти можуть бути наслідком помилок у даних, шахрайських дій, несправностей у технічних

системах, або ж сигналізувати про нові явища, які ще не мають пояснення, але можуть бути критично важливими для бізнесу чи безпеки. Наприклад, одинична підозріла транзакція у банківській системі може виявитися спробою зламу, а раптове зростання температури у вузлі обладнання — ознакою аварійної ситуації [9].

Обидва ці підходи — кластеризація та виявлення аномалій — зазвичай застосовуються не ізольовано, а як частина більшого процесу аналізу даних, який включає попередню підготовку вибірки. Зокрема, для підвищення ефективності моделей часто застосовується нормалізація числових ознак, усунення пропущених значень, видалення викидів, відбір релевантних характеристик або зниження розмірності за допомогою таких методів, як Principal Component Analysis (PCA) [4]. Коректна обробка даних є запорукою надійності результатів, особливо в умовах реальних, "брудних" даних, де поширені неоднорідність, шум, пропуски й надлишковість [3]. Для цього використовуються стандартні засоби аналізу даних, зокрема бібліотеки Python, які будуть детально розглянуті у наступних розділах.

1.2. Кластеризація

Кластеризація (англ. *clustering*) — це метод аналізу даних, що дозволяє виявляти структуру в наборі немаркованих об'єктів шляхом групування їх за певними критеріями подібності. На відміну від класифікації, де кожен об'єкт має заздалегідь відомий клас, кластеризація не потребує попередньої інформації про мітки даних. [6]

Метою кластеризації є поділ множини об'єктів $X = \{x_1, x_2, \dots, x_n\}$ на підмножини C_1, C_2, \dots, C_k , де кожен кластер C_i складається з об'єктів, подібних між собою згідно з вибраною метрикою (наприклад, евклідова відстань).

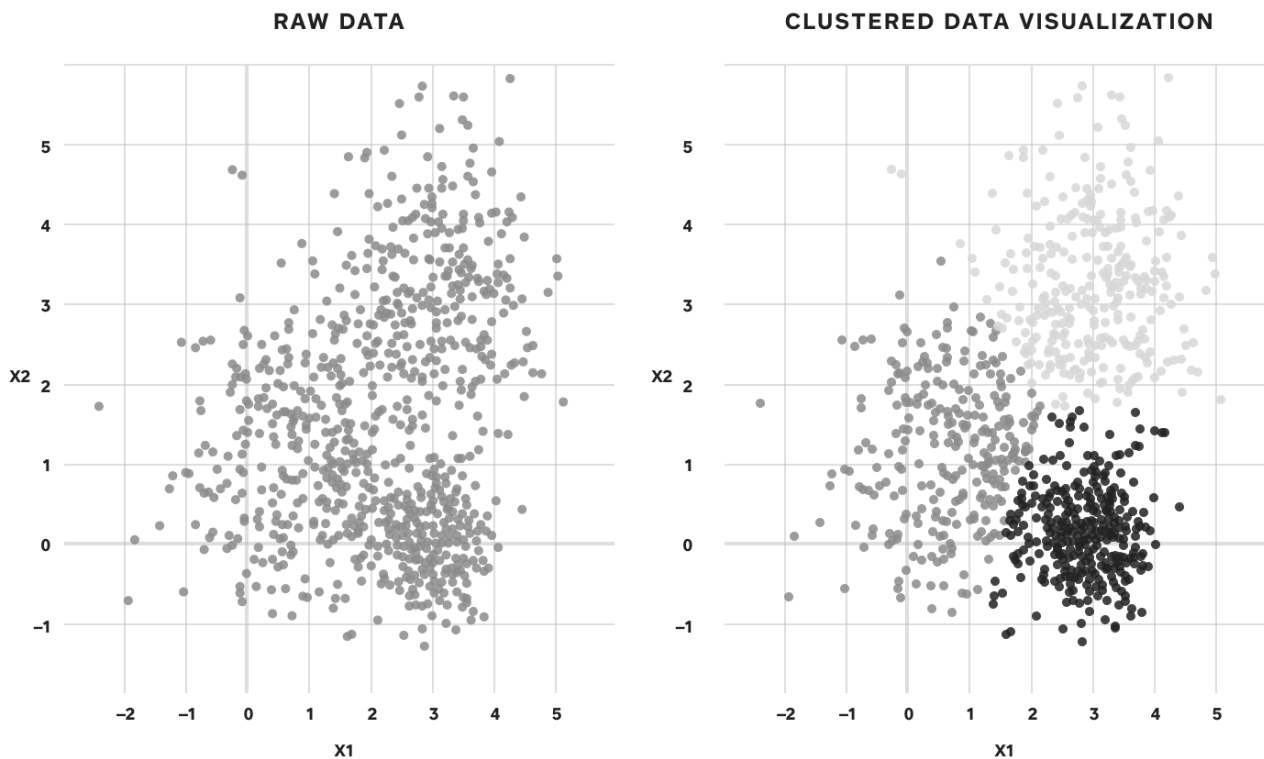


Рисунок 1.1 – Приклад візуалізації даних до і після кластеризації

Типи кластерів

- Компактні кластери: об'єкти згруповані навколо центру (наприклад, у K-Means).
- Щільні кластери: об'єкти згруповані в області з високою густиною (наприклад, у DBSCAN).
- Ієрархічні структури: багаторівневе групування з побудовою дендрограми.

Основні алгоритми кластеризації

K-Means

Алгоритм **K-Means** є одним із найвідоміших і найбільш застосовуваних методів кластерного аналізу. Його метою є розподіл множини об'єктів на наперед задану кількість кластерів таким чином, щоб об'єкти в межах одного кластеру були подібними, а між різними кластерами — відмінними за обраною метрикою. Найчастіше використовується евклідова відстань [4].

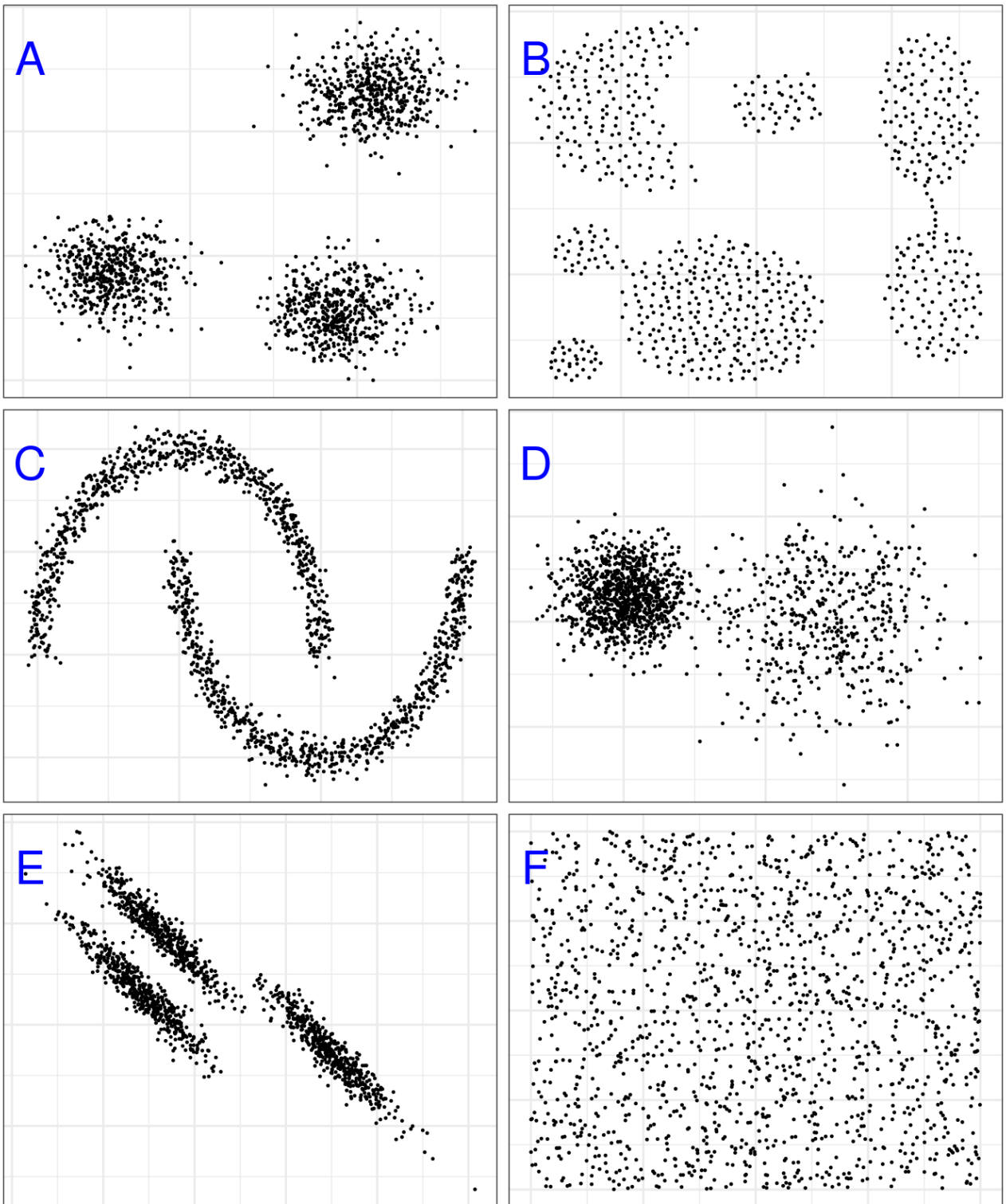


Рисунок 1.2 — Приклади кластерних структур у двовимірному просторі (А, краплі; В, агрегація; С, шумні супутники; D, різна щільність; E, анізотропні розподіли; F, відсутність структури.)

Процес кластеризації в K-Means включає випадкову ініціалізацію центрів кластерів, подальше призначення кожного об'єкта до найближчого центру, обчислення нових центрів як середніх значень ітераційно до досягнення збіжності. Зазвичай алгоритм сходиться швидко, проте його результат залежить від початкового розміщення центрів, що може призводити до локальних мінімумів [7].

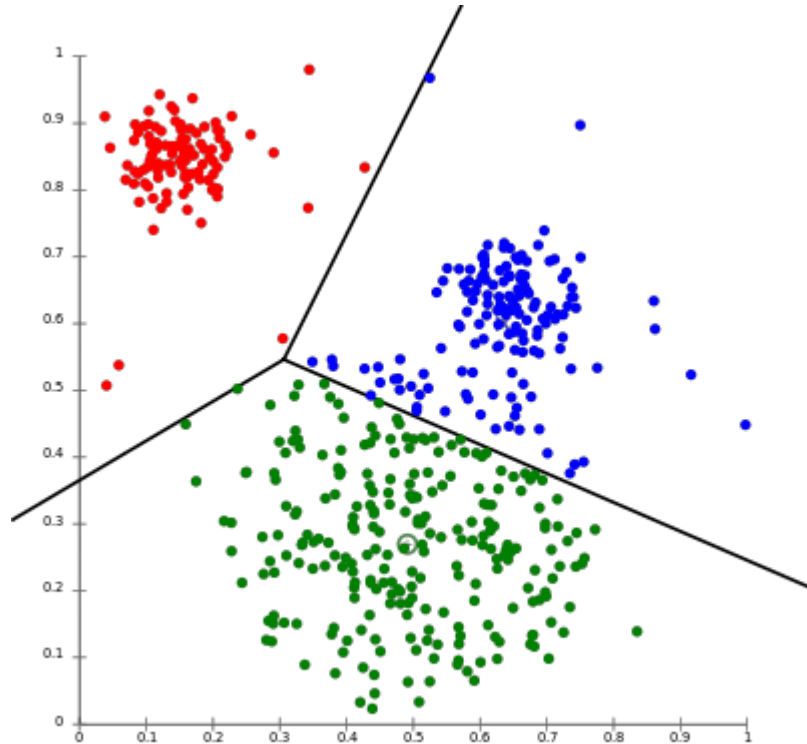


Рисунок 1.3 — Результат кластеризації K-Means

K-Means добре працює з опуклими, однорідними кластерами, але має обмеження при роботі з кластерами складної форми або неоднорідної густини [6].

$$\sum_{i=1}^k \sum_{x_j \in C_i} \|x_j - \mu_i\|^2 \quad (1.1)$$

де μ_i – центр кластера C_i

DBSCAN (Density-Based Spatial Clustering of Applications with Noise)

Алгоритм **DBSCAN** (*Density-Based Spatial Clustering of Applications with Noise*) належить до класу методів кластеризації, що ґрунтуються на щільності. Його ключова ідея полягає в тому, що кластери — це області простору з високою густиною об'єктів, відокремлені зонами низької щільності. На відміну від K-Means, DBSCAN не потребує задання кількості кластерів заздалегідь та дозволяє виявляти кластери довільної форми.

Алгоритм визначає дві ключові величини: ϵ — радіус околу, в межах якого перевіряється щільність, та **MinPts** — мінімальну кількість точок, необхідних для формування кластера. Точки, що не входять до жодного кластера, позначаються як шум. DBSCAN стійкий до викидів і добре працює в задачах, де кластери мають різні форми та розміри, однак може бути чутливим до вибору параметрів [7].

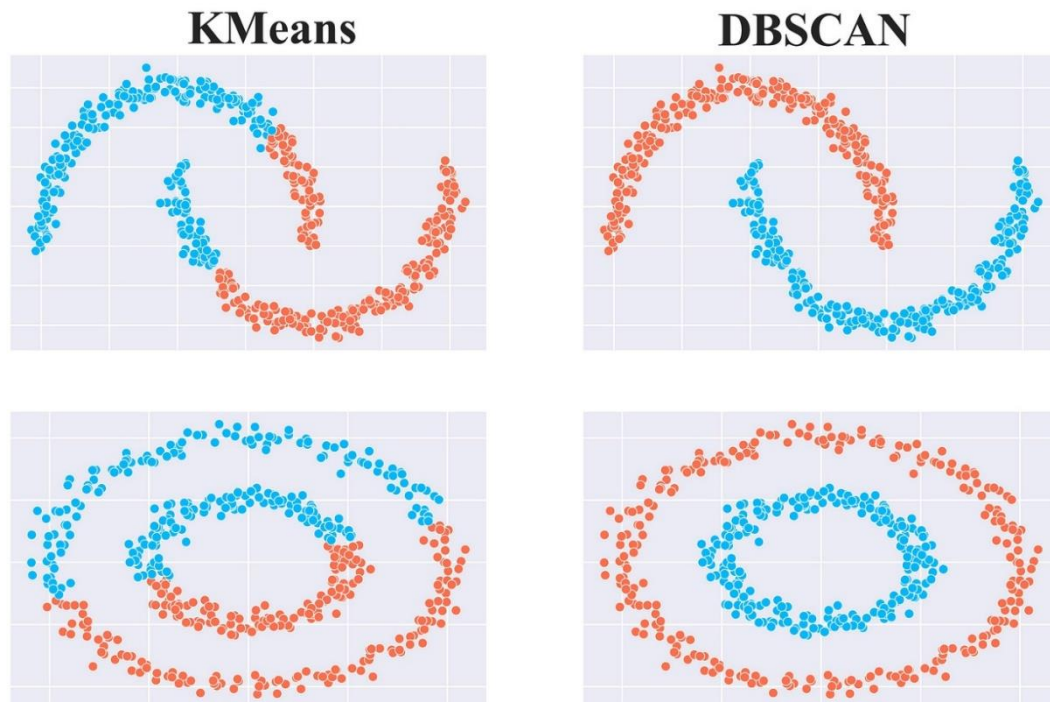


Рисунок 1.4 — Порівняння кластеризації методами K-Means та DBSCAN

Ієрархічна кластеризація

Ієрархічна кластеризація — це підхід до групування об'єктів, який не потребує попереднього задання кількості кластерів і буде вкладену структуру

кластерів у вигляді дерева. Така структура називається **дендрограмою** і дозволяє досліднику візуально обрати бажаний рівень кластеризації, залежно від глибини об'єднання.

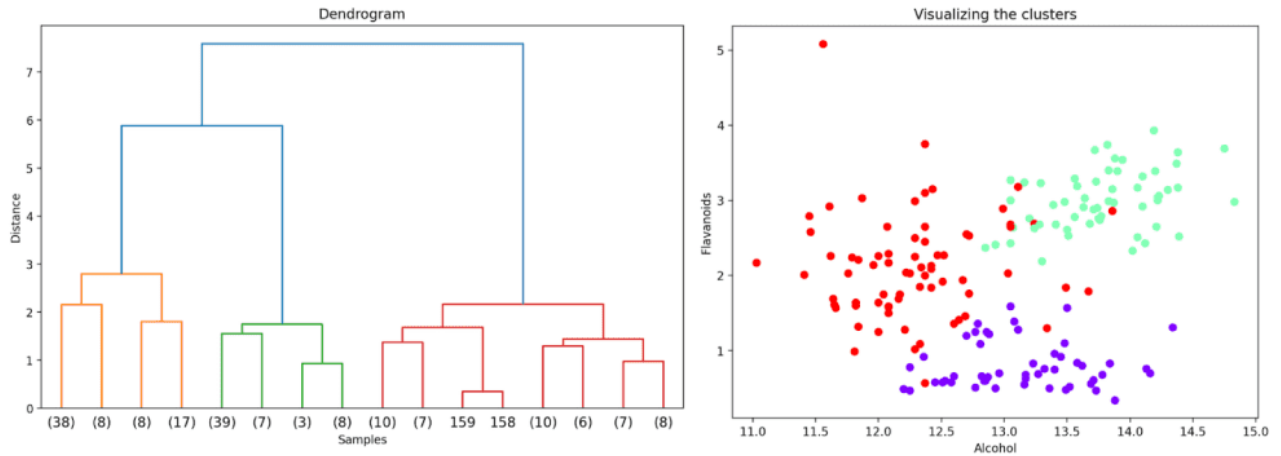


Рисунок 1.5 — Візуалізація ієрархічної кластеризації

Існують два основні варіанти ієрархічного аналізу: **агломеративний** (знизу догори) та **дивізивний** (згори донизу). У першому випадку кожен об'єкт спочатку є окремим кластером, які поступово об'єднуються; у другому — всі об'єкти спочатку розглядаються як один кластер, який послідовно розділяється. Як міру відстані між кластерами використовують середню, мінімальну або максимальну відстань між їхніми елементами. Метод добре працює для невеликих наборів даних, але має обмежену масштабованість через обчислювальну складність [8].

Підходи до вибору кількості кластерів

- Метод **ліктя** (elbow method): аналіз графіка залежності дисперсії від числа кластерів.
- Силуетний метод (silhouette score): оцінка якості кластеризації.
- Крос-валідація за допомогою класифікатора (у спеціалізованих задачах).

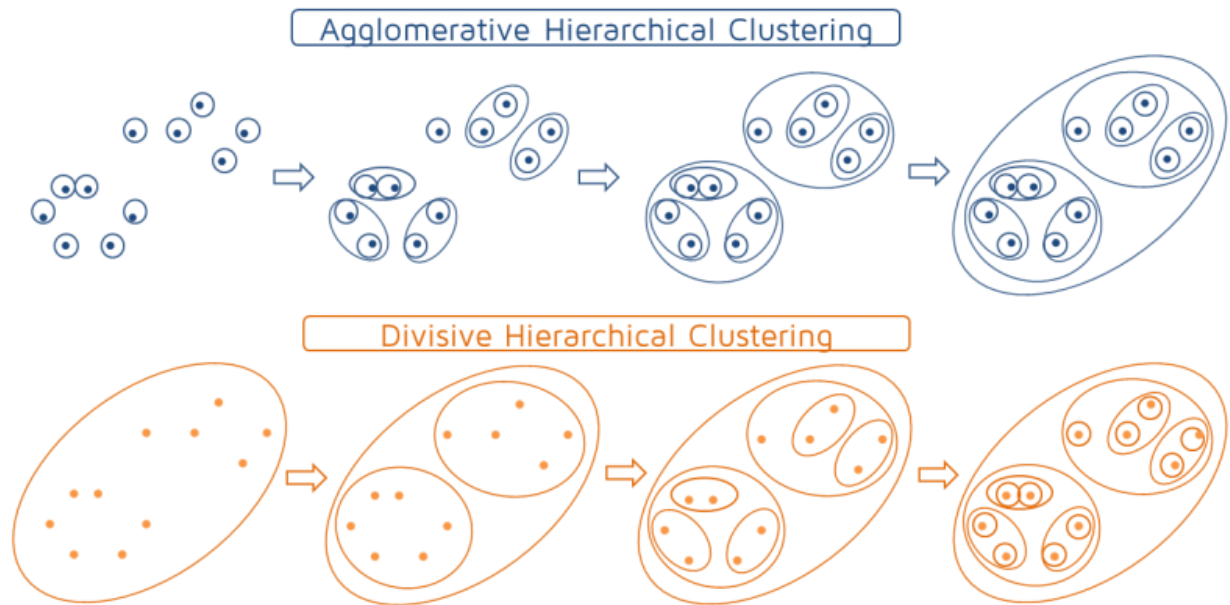


Рисунок 1.6 — Основні варіанти ієрархічного аналізу

1.3. Виявлення аномалій

Виявлення аномалій (англ. *anomaly detection* або *outlier detection*) — це задача виявлення об'єктів або спостережень, які істотно відрізняються від інших за певними характеристиками. Такі об'єкти можуть бути наслідком помилок у даних, але також часто є джерелом критично важливої інформації: ознакою зламу системи, фінансового шахрайства, раптової відмови технічного пристрою або появи нової тенденції.

Залежно від контексту, аномалії можуть мати різну природу:

- **Точкові аномалії** — окремі об'єкти суттєво відрізняються від решти. Наприклад, один надто великий платіж серед багатьох звичайних.
- **Контекстні аномалії** — об'єкти є аномальними в конкретному контексті, наприклад, температура 30°C — нормальна для літа, але аномальна для зими.

- **Колективні аномалії** — група об'єктів разом утворює аномальну поведінку, хоча кожен окремо може виглядати нормально (наприклад, підозрілий шаблон мережевого трафіку).

Методи виявлення аномалій

Існує кілька класів підходів до виявлення аномалій, кожен з яких має свої переваги та обмеження.

Статистичні методи

Базуються на припущенні, що нормальні дані відповідають деякому ймовірнісному розподілу. Якщо точка має низьку ймовірність, вона вважається аномальною.

- **Z-score**: точка вважається аномалією, якщо її відхилення від середнього перевищує певний поріг:

$$z = \frac{x - \mu}{\sigma} \quad (1.1)$$

- **IQR (interquartile range)**: використовується для виявлення викидів у розподілах.

Переваги: Простота та інтерпретованість

Недоліки: Неefективність у складних, багатовимірних розподілах

Кластерні методи

Базуються на припущенні, що нормальні об'єкти лежать всередині щільних кластерів, а аномалії — далеко від них.

- **K-Means distance-based anomaly score**: якщо точка знаходиться далеко від центру свого кластера, вона вважається підозрілою.

- **DBSCAN**: точки, які не належать до жодного кластеру, можна вважати аномальними.

Моделі машинного навчання

Ці методи зазвичай працюють без учителя (unsupervised), але можливе і навчання з учителем (у рідкісних випадках, коли є мітки).

- **Isolation Forest**: будує випадкові дерева для “ізоляції” об’єктів. Аномалії ізолюються за меншу кількість кроків.
- **Local Outlier Factor (LOF)**: оцінює щільність навколишнього простору точки й порівнює її з іншими. Якщо точка ізольована, вона вважається аномальною.

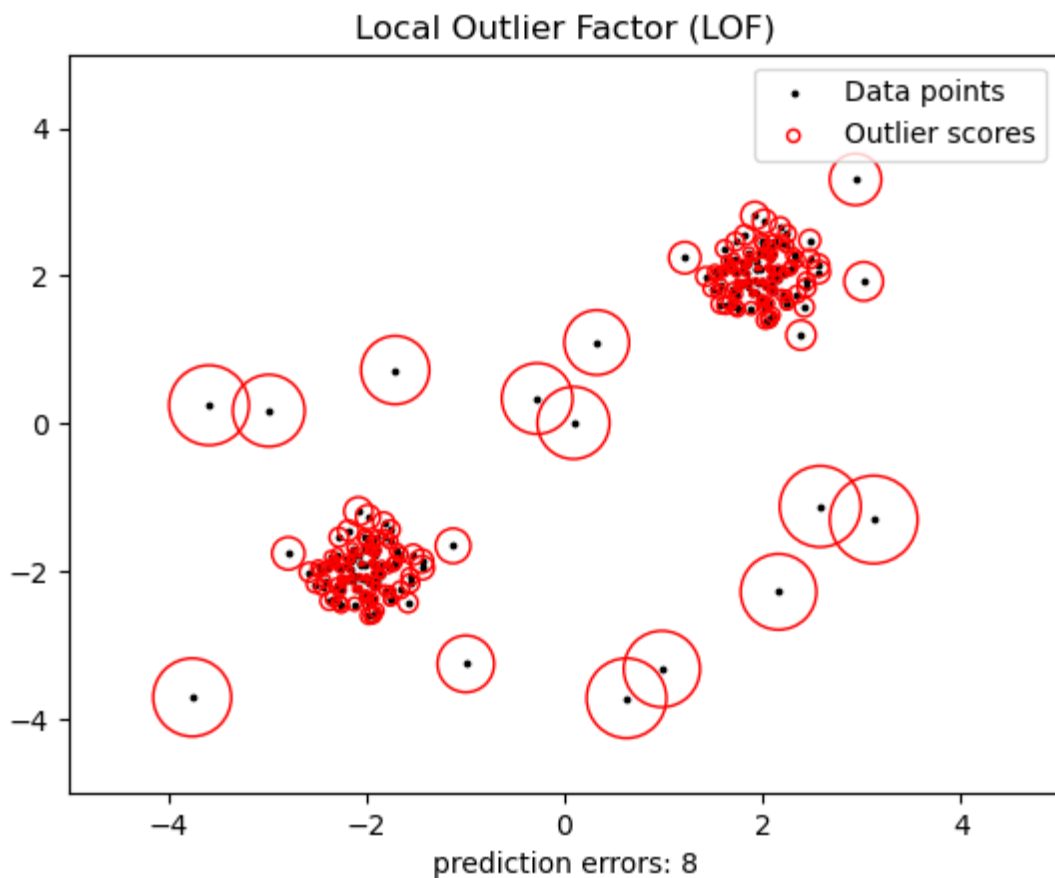


Рисунок 1.7 — Приклад виявлення аномалій методом LOF

Підходи на основі нейронних мереж

У задачах з великою кількістю ознак або складними зв'язками між змінними використовуються автокодера (autoencoders), які навчаються відновлювати вхідні дані. Якщо модель не може точно відновити дані — вони можуть бути аномальними.

Застосування:

- Медицина (виявлення рідкісних патологій)
- Кібербезпека (вторгнення в систему)
- Технічна діагностика (виявлення поломок)

Проблеми і виклики

- Наявність високовимірних даних (curse of dimensionality)
- Баланс між виявленням аномалій і хибнопозитивними результатами
- Відсутність еталонних аномалій у тренувальних вибірках

1.4. Огляд існуючих рішень

Сучасні інструменти для аналізу даних пропонують широкий спектр методів кластеризації та виявлення аномалій. Існують як повноцінні програмні платформи з графічним інтерфейсом, так і гнучкі бібліотеки для програмної реалізації в середовищах розробки. У цьому підрозділі розглянемо найбільш популярні з них, особливо ті, що можуть бути використані як основа або джерело ідей для створення власного застосунку.

Scikit-learn (Python)

Бібліотека Scikit-learn є одним із найпоширеніших інструментів для машинного навчання в мові Python. Вона реалізує широкий набір алгоритмів, включно з кластеризацією та виявленням аномалій.

Можливості:

- KMeans, MiniBatchKMeans, DBSCAN, Agglomerative Clustering

- IsolationForest, LocalOutlierFactor
- Підтримка Pipelines, GridSearch для налаштування параметрів
- Підтримка візуалізації результатів (разом із matplotlib або seaborn)

Переваги:

- Простий у використанні API
- Широке ком'юніті
- Документація з прикладами

ELKI (Java)

ELKI (Environment for Developing KDD-Applications Supported by Index-Structures) — це академічна Java-платформа, орієнтована на кластеризацію, виявлення аномалій та індексні структури.

Переваги:

- Великий набір новітніх алгоритмів, яких немає в інших інструментах
- Підтримка параметричного дослідження

Недоліки:

- Високий поріг входу
- Не призначений для широкого кола користувачів

Orange

Orange — візуальне середовище аналізу даних, засноване на Python.

Переваги:

- Інтуїтивно зрозумілий графічний інтерфейс
- Підтримка кластеризації, PCA, візуалізацій

Недоліки:

- Менша гнучкість у порівнянні з чистим Python API
- Обмежені можливості для кастомізації логіки

RapidMiner

RapidMiner — платформа для аналітики, орієнтована на корпоративних користувачів. Підтримує drag-and-drop моделювання, має модулі для кластеризації, виявлення аномалій, прогнозування.

Недоліки:

- Комерційна ліцензія для повного функціоналу
- Низька придатність для вбудовування у кастомні рішення

PyOD (Python Outlier Detection)

PyOD — спеціалізована бібліотека Python для виявлення аномалій. Містить понад 30 методів: kNN, LOF, COPOD, AutoEncoder, NBOS, ECOD та інші.

Переваги:

- Широке покриття методів аномалій
- API схожий на scikit-learn
- Підтримка комбінованих моделей (ensemble)

У таблиці 1.1 наведено порівняння основних характеристик зазначених інструментів:

Таблиця 1.1 - Порівняння засобів для виявлення аномалій

Інструмент	Кількість алгоритмів	Зручність використання	Підтримка візуалізації	Якість документації
Scikit-learn	20+ кластеризація, 2–3 для аномалій	Висока: уніфіковане API, підтримка Pipelines	Часткова (через matplotlib, seaborn)	Висока: туторіали, API, приклади [11]
PyOD	30+ алгоритмів виявлення аномалій	Висока: сумісний зі Scikit-learn	Обмежена (через власні методи або matplotlib)	Висока: чітка структура, приклади [12]
ELKI	100+ алгоритмів	Середня: потребує знань Java, конфігурації через CLI	Обмежена: зовнішній вивід, без інтерактиву	Середня: орієнтація на дослідників [5]

Інструмент	Кількість алгоритмів	Зручність використання	Підтримка візуалізації	Якість документації
Orange	~10 основних, розширення через add-ons	Висока: графічний інтерфейс, drag-and-drop	Висока: інтерактивні діаграми «з коробки»	Середня: базова документація, community [3]
RapidMiner	~20 в базовій версії, більше у розширеннях	Висока: візуальне моделювання, блокова логіка	Висока: інтегровані графіки, auto-visualizations	Середня: частково обмежена в безкоштовній версії [14]

Як видно з таблиці, для розробки прикладного програмного забезпечення на мові Python найбільш доцільними є Scikit-learn і PyOD. Вони забезпечують широкий вибір алгоритмів, підтримують інтеграцію один з одним, мають прозоре API та активну спільноту. Інструменти на зразок Orange та RapidMiner більше підходять для візуального моделювання й навчальних цілей, у той час як ELKI залишається потужним, проте складним фреймворком, переважно орієнтованим на наукові дослідження.

У розділі проведено огляд основних теоретичних підходів та практичних рішень, що стосуються кластеризації та виявлення аномалій — двох важливих задач аналізу даних, які мають широке застосування в сучасних інформаційних системах.

Під час вивчення кластеризації проаналізовані основні алгоритми, зокрема K-Means, DBSCAN та ієрархічна кластеризація. Кожен з них має свої переваги та недоліки, які визначають доцільність використання залежно від типу даних, їх розмірності, структури та очікуваних результатів. K-Means підходить для компактних кластерів з приблизно рівною густиною, DBSCAN — для кластерів довільної форми і виявлення шуму, а ієрархічні методи забезпечують багаторівневе уявлення про структуру даних.

Аналіз методів виявлення аномалій показав, що ефективність різних підходів залежить від природи аномалій та розподілу даних. Статистичні методи

добре працюють на одновимірних даних із нормальною структурою, тоді як більш гнучкі методи, такі як Local Outlier Factor, Isolation Forest або автокодери, краще підходять для складних багатовимірних вибірок. Особливої уваги потребує вибір метрик подібності та способів оцінки "нормальності" об'єкта.

Огляд існуючих рішень показав, що найбільш доцільними з точки зору реалізації власного застосунку є бібліотеки `scikit-learn` (для кластеризації) та `PyOD` (для виявлення аномалій). Ці інструменти реалізовані мовою Python, мають зручний API, сумісні між собою та підтримують інтеграцію з іншими бібліотеками Python-екосистеми. Крім того, вони активно підтримуються спільнотою та мають достатню продуктивність для задач, які не вимагають обробки великих обсягів даних у реальному часі.

2 ОБҐРУНТУВАННЯ ВИБОРУ ТЕХНОЛОГІЙ І ІНСТРУМЕНТІВ

2.1. Огляд мов програмування для аналізу даних

У сучасній практиці розробки прикладного програмного забезпечення для обробки й аналізу даних надзвичайно важливим є вибір мови програмування. Особливо це актуально в задачах кластеризації та виявлення аномалій, де потрібна не лише обчислювальна ефективність, а й наявність готових бібліотек, зручних засобів для роботи з даними, побудови візуалізацій та швидкої реалізації прототипів. Мова програмування визначає не лише синтаксис і стиль коду, а й доступність рішень, інструментів, продуктивність, можливість інтеграції з іншими системами. У цьому підрозділі буде здійснено поетапний огляд чотирьох мов програмування — Python, R, Java та C++ — які є найпопулярнішими у сфері Data Science, та визначено, яка з них найкраще відповідає поставленій задачі.

Python

Python є однією з найпопулярніших мов програмування у сфері аналізу даних, машинного навчання та прикладного штучного інтелекту. Вона створена як мова загального призначення, але з часом набула широкого поширення в науковій спільноті та в індустрії завдяки простому синтаксису, гнучкості та активному розвитку відкритих бібліотек.

Однією з головних переваг Python є його бібліотечна екосистема. Для задач кластеризації та виявлення аномалій доступні численні пакети: scikit-learn, PyOD, hdbscan, numpy, pandas, matplotlib, seaborn, plotly та інші. Наприклад, scikit-learn реалізує алгоритми K-Means, DBSCAN, агломеративну кластеризацію, а також виявлення аномалій через Local Outlier Factor та Isolation Forest [4]. Бібліотека PyOD містить понад 30 спеціалізованих методів для виявлення викидів у табличних даних, включаючи класифікатори на основі глибокого навчання.

Ще однією перевагою є **зрозумілий синтаксис**, який робить Python доступним навіть для новачків. Завдяки цьому створення прототипів систем, побудова тестів, моделювання та аналіз відбуваються значно швидше, ніж у більшості інших мов. Код, написаний на Python, зазвичай коротший та легший для читання, що спрощує підтримку програмного забезпечення.

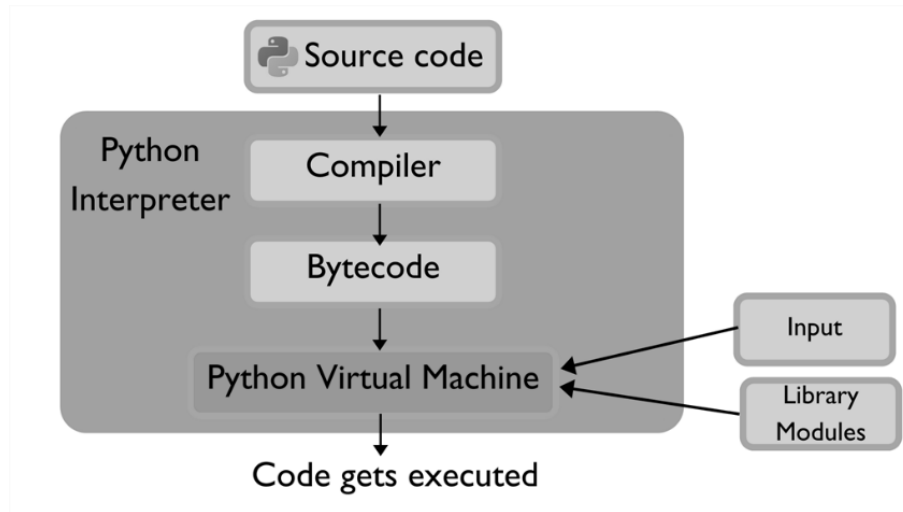


Рисунок 2.1 Структурна схема функціонування мови Python

У Python реалізовано широкий набір засобів візуалізації: `matplotlib` для базових графіків, `seaborn` для статистичних діаграм, `plotly` для інтерактивних дашбордів, `bokeh` для веб-візуалізацій. Це дозволяє швидко аналізувати кластери, виявлені аномалії та структуру даних.

Python має активну глобальну спільноту, регулярно оновлювану документацію, а також велику кількість онлайн-курсів, форумів, відкритих наборів даних і прикладів використання. Це дозволяє суттєво скоротити час пошуку рішень на типові проблеми.

Щодо продуктивності, Python не є найшвидшою мовою виконання, але багато критичних бібліотек написано на C/C++, що забезпечує достатню ефективність. Для задач, де продуктивність критична, використовуються гібридні підходи — наприклад, Python як обгортка до високопродуктивного коду або використання бібліотек із GPU-підтримкою.

У цілому, Python поєднує простоту синтаксису, потужність бібліотек, універсальність і швидкість розробки, що робить його найбільш придатним для реалізації системи кластеризації та виявлення аномалій.

R

Мова R була розроблена для статистичних обчислень та графічного аналізу і має глибоке коріння в академічному середовищі. Вона активно використовується у дослідницьких установах, економічному аналізі, біостатистиці, психології та соціології.

R має потужну базу пакетів для кластеризації та аномалій, зокрема `cluster`, `factoextra`, `anomalize`, `TSclust`, `fpc`, `dbscan`. У цих бібліотеках реалізовано алгоритми K-Means, DBSCAN, hierarchical clustering, LOF та інші. Проте, через специфічну орієнтацію на статистику, бібліотеки часто вимагають глибшого розуміння статистичних методів і мають складніший API для інженерів [3].

Окремо варто відзначити візуалізаційні можливості R. `ggplot2` є однією з найсильніших бібліотек у світі для побудови публікаційних графіків. Візуалізація в R є більш гнучкою для статистичних експериментів, однак менш зручною в умовах побудови програмних інтерфейсів або інтерактивних систем.

R поступається Python у контексті побудови масштабованих застосунків. Мова не призначена для розробки повноцінних систем або взаємодії з API, базами даних, вебінтерфейсами — хоча такі можливості й існують, вони реалізовані складніше.

Також у R значно менше інженерних інструментів, таких як менеджери залежностей, віртуальні середовища, інструменти CI/CD. Це створює труднощі в командній розробці, особливо у складних проєктах.

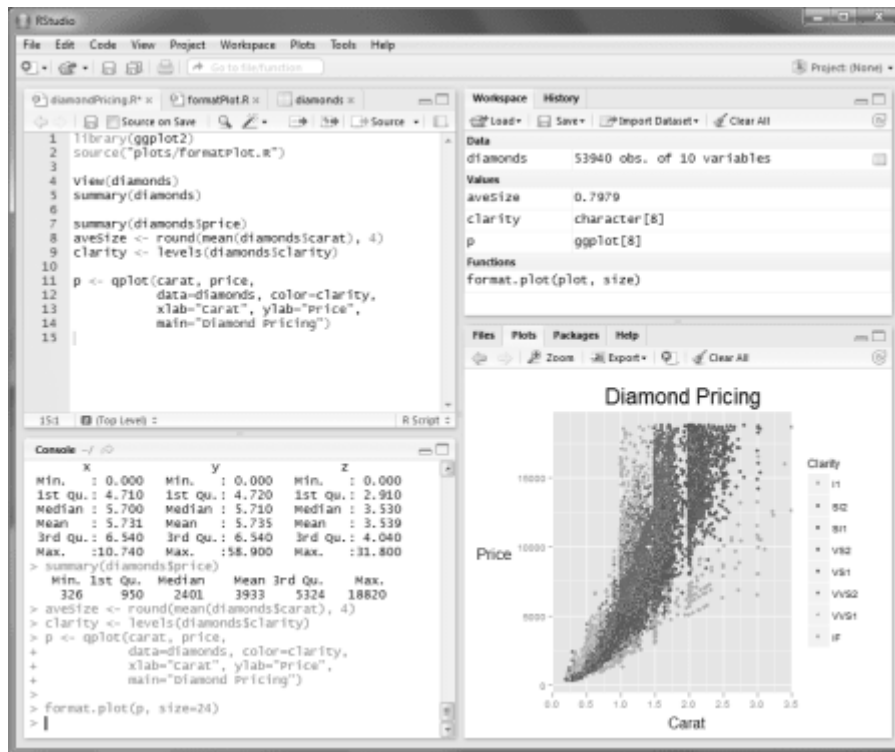


Рисунок 2.2 Середовище розробки R-Studio

Підсумовуючи, R є ефективним інструментом для статистичного аналізу, швидкої побудови графіків та обробки невеликих наборів даних, але менш придатна для повноцінного програмного продукту або інженерного рішення.

Java

Java — одна з найпоширеніших об'єктно-орієнтованих мов програмування, яку активно використовують у великих корпоративних застосунках, банківських системах, серверних рішеннях та мобільній розробці. Вона вирізняється високою продуктивністю, стабільністю, переносимістю та багаторівневою підтримкою безпеки.

У сфері аналізу даних Java представлена кількома популярними бібліотеками та фреймворками, такими як **Weka**, **ELKI**, **Deeplearning4j**, **Apache Mahout**. Зокрема, **Weka** надає графічний інтерфейс і реалізує багато класичних алгоритмів кластеризації: K-Means, Hierarchical Clustering, DBSCAN. **ELKI** є високоспеціалізованим фреймворком, розробленим для науковців, і включає

десятки алгоритмів, яких немає в інших мовах, включно з великою кількістю методів для виявлення аномалій [5].

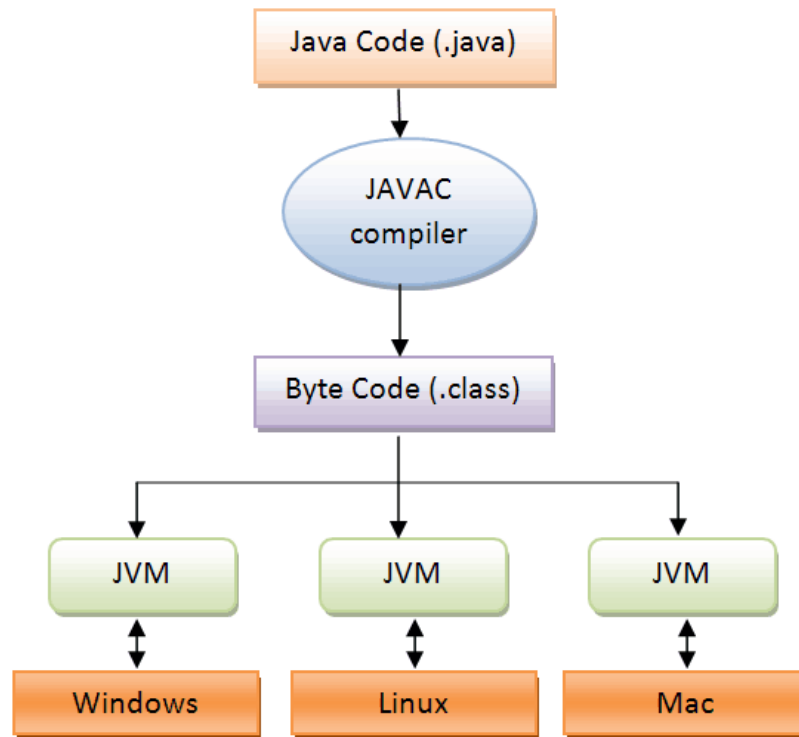


Рисунок 2.3 Схема функціонування мови Java

Однак Java має низку суттєвих **обмежень** для задач швидкої побудови прототипів. Синтаксис є об'ємним і вимагає значної кількості шаблонного коду. Написання навіть базових операцій — як зчитування даних чи побудова графіків — займає більше часу, ніж у Python чи R. Крім того, в Java бракує зручних нативних засобів для **візуалізації результатів кластеризації**, тому часто доводиться або експортувати дані у CSV для подальшого аналізу в іншому середовищі, або підключати зовнішні інструменти.

Хоча Java є ефективною з погляду **виконання** та добре працює з багатопоточністю, її **використання для прикладної аналітики ускладнене**. Навіть Weka або ELKI, попри свою функціональність, мають вузьку

спеціалізацію й недостатньо розвинені засоби інтеграції з іншими сучасними ML-системами.

Отже, Java може використовуватись для реалізації **окремих вузлів складних систем**, проте вона не є найкращим вибором для швидкої розробки інструменту кластеризації з виводом результатів, гнучкою обробкою даних і зручним інтерфейсом.

C++

C++ — це мова системного програмування, яка забезпечує найвищу продуктивність серед усіх розглянутих варіантів. Вона дозволяє працювати з пам'яттю на низькому рівні, контролювати потоки виконання та ефективно реалізовувати алгоритми. Саме тому C++ використовується в комп'ютерному зорі, обробці сигналів, ігровій індустрії, системах із жорсткими обмеженнями по ресурсах.

У контексті аналізу даних, C++ має обмежену популярність, однак існують кілька бібліотек, які дозволяють реалізувати кластеризацію та пошук аномалій. Зокрема, можна згадати **mlpack**, **dlib**, **Shark** та **OpenCV** (частково). Наприклад, **mlpack** реалізує K-Means, DBSCAN та інші базові алгоритми, забезпечуючи високу швидкодію. Однак використання цих бібліотек вимагає складного налаштування, вміння працювати з шаблонами та компіляторами, а також поглибленого розуміння внутрішньої структури алгоритмів [6].

Основною проблемою C++ у контексті кластерного аналізу є **висока вартість розробки**: створення простого застосунку з інтерфейсом командного рядка, побудовою графіків, роботою з CSV-файлами та реалізацією алгоритмів вимагає значно більше зусиль, ніж у високорівневих мовах.

C++ практично не має **інтегрованих інструментів візуалізації**, що змушує розробників експортувати дані для обробки в Python або R, або реалізовувати власні графічні модулі. Це знижує гнучкість і прискорює накопичення технічного боргу в проєктах.

Ще одним недоліком є **низька доступність фреймворків для виявлення аномалій**. Хоча можливо реалізувати потрібні алгоритми вручну, це потребує глибокої математичної підготовки та великого обсягу коду.

Таким чином, C++ — це **ефективний, але надмірно складний** варіант для задач кластеризації в контексті прикладної розробки. Його застосування доцільне лише у випадках, коли критичною є максимальна продуктивність.

Таблиця 2.1 – Порівняння мов програмування

Критерій	Python	R	Java	C++
Наявність бібліотек	Висока	Висока	Середня	Низька
Робота з табличними даними	Зручна	Зручна	Ускладнена	Ускладнена
Візуалізація	Потужна	Дуже потужна	Слабка	Відсутня
Складність синтаксису	Низька	Середня	Висока	Висока
Продуктивність	Середня	Середня	Висока	Дуже висока
Швидкість розробки	Висока	Середня	Низька	Дуже низька
Підтримка спільноти	Велика	Обмежена	Помірна	Низька
Підтримка аномалій з коробки	Так	Частково	Обмежено	Ні
Гнучкість і розширюваність	Висока	Середня	Низька	Низька
Інструменти автоматизації	Добре розвинені	Слабо розвинені	Добре	Слабо

Проведений аналіз свідчить, що кожна з розглянутих мов програмування має свої переваги та сфери застосування. R добре підходить для глибокого статистичного аналізу та візуалізацій, Java — для розробки великих корпоративних систем, C++ — для критичних до продуктивності програм.

Однак саме Python забезпечує оптимальний баланс між простотою використання, швидкістю розробки, підтримкою алгоритмів кластеризації та

виявлення аномалій, бібліотеками для візуалізації, обробки даних, а також широкою підтримкою спільноти.

Саме тому в межах даної кваліфікаційної роботи обрано мову Python, як таку, що найповніше відповідає вимогам проєкту та забезпечує можливість реалізації функціонального, ефективного й масштабованого програмного застосунку для аналізу даних.

2.2. Бібліотеки та програмні засоби для реалізації застосунку

У сучасному програмуванні одним із головних факторів, що визначають швидкість та якість розробки прикладного програмного забезпечення, є доступність готових бібліотек та інструментів. Особливо це важливо у сфері аналізу даних, де ключову роль відіграють алгоритми машинного навчання, засоби обробки структурованої інформації та візуалізації.

Мова програмування Python, яка була обґрунтовано обрана в попередньому підрозділі, вирізняється потужною бібліотечною екосистемою, яка дозволяє реалізовувати складні аналітичні задачі з мінімальними витратами часу. У цьому підрозділі буде розглянуто основні бібліотеки, які використовуються для реалізації кластеризації, виявлення аномалій, попередньої обробки даних та візуалізації результатів. Окремо буде описано їхні функціональні можливості, переваги, специфіку використання та причини включення до структури розроблюваного застосунку.

scikit-learn

scikit-learn — одна з найпопулярніших бібліотек Python для реалізації алгоритмів машинного навчання. Вона створена на основі таких базових бібліотек, як NumPy, SciPy та matplotlib, і пропонує зручний інтерфейс для широкого спектра алгоритмів навчання без учителя (unsupervised learning), у тому числі кластеризації.

Бібліотека включає реалізації таких алгоритмів, як:

- K-Means, включаючи MiniBatchKMeans для великих даних;
- DBSCAN — кластеризація за щільністю;
- Agglomerative Clustering — агломеративна ієрархічна кластеризація;
- Spectral Clustering — спектральні методи кластеризації;
- Birch — ефективна кластеризація великих обсягів даних.

Окрім кластеризації, scikit-learn надає доступ до кількох алгоритмів виявлення аномалій, зокрема:

- Isolation Forest — метод ізоляції викидів;
- Local Outlier Factor (LOF) — визначає ступінь ізольованості точки;
- One-Class SVM — метод на основі підтримуючих векторів для визначення області нормальних даних.

Бібліотека підтримує численні утиліти для попередньої обробки, зокрема нормалізацію, масштабування, заповнення пропущених значень, а також забезпечує модулі для перехресної валідації, оцінки якості моделей, побудови pipelines.

scikit-learn підтримує єдиний уніфікований API, що значно спрощує експериментування, зміну параметрів і заміну алгоритмів без переписування коду. Вона добре документована, має великий набір прикладів і активно підтримується спільнотою [11].

PyOD

PyOD (Python Outlier Detection) — спеціалізована бібліотека для виявлення аномалій. Вона містить понад 30 алгоритмів, включаючи як класичні статистичні методи, так і сучасні підходи з використанням ансамблів і глибокого навчання.

Найважливішими реалізованими методами є:

- KNN, LOF, ABOD, Feature Bagging — класичні методи;
- AutoEncoder, VAE (Variational AutoEncoder) — глибокі нейромережеві архітектури;

- ECOD, COPOD, HBOS, IForest, SOD, MCD тощо.

PyOD є сумісною з scikit-learn, використовує той самий інтерфейс `.fit()`, `.predict()`, `.decision_function()`, що дозволяє легко інтегрувати її в існуючі пайплайни.

Особливою перевагою PyOD є вбудовані засоби валідації: вона дозволяє легко оцінювати ефективність методів на синтетичних або реальних наборах даних, використовуючи метрики `precision`, `recall`, `AUC` тощо. Також підтримується візуалізація результатів виявлення аномалій, що є корисним для аналізу якості моделі.

Завдяки зручному дизайну, модульності та великому вибору методів, PyOD є ключовим інструментом для задач аномалій у цьому проєкті [12].

pandas

`pandas` — базова бібліотека для роботи з табличними даними в Python. Вона забезпечує структуру `DataFrame`, яка дозволяє зручно маніпулювати наборами даних, як у реляційних СУБД або Excel.

Основні можливості:

- імпорт даних з форматів CSV, Excel, JSON, SQL;
- фільтрація, агрегація, сортування, трансформація таблиць;
- обробка пропущених значень (заповнення, видалення, інтерполяція);
- побудова нових ознак (`feature engineering`);
- групування за категоріями (`groupby`);
- інтеграція з scikit-learn через `values` або `.to_numpy()`.

`pandas` забезпечує високу гнучкість у підготовці даних: саме вона дозволяє очищувати вибірки перед кластеризацією та масштабувати їх для подальшого аналізу. Без цього етапу використання алгоритмів машинного навчання неефективне або навіть некоректне.

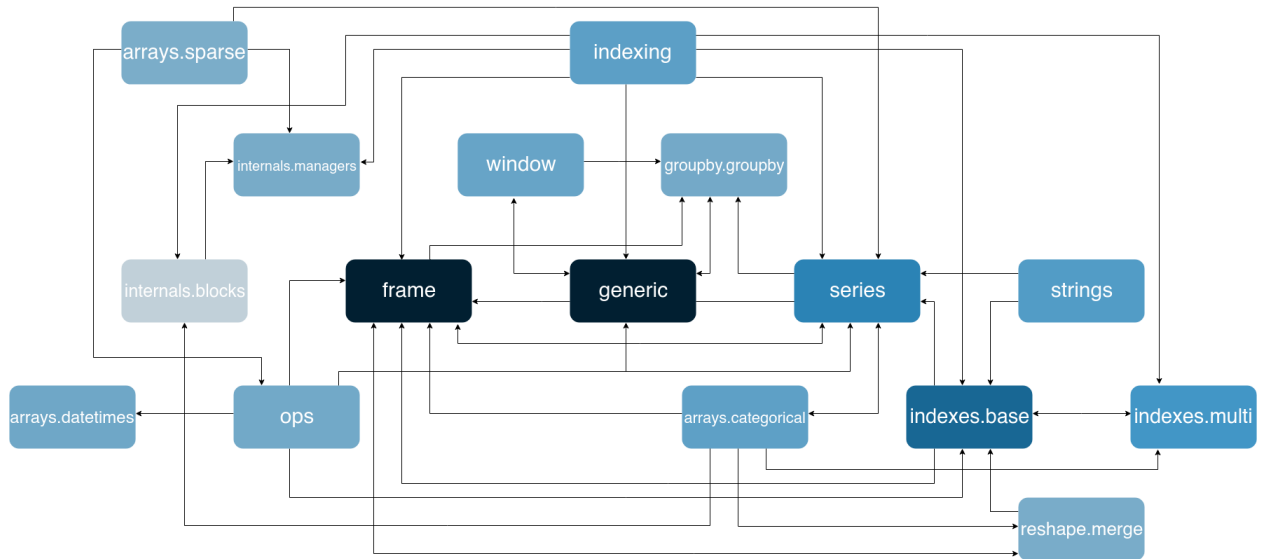


Рисунок 2.4 – Взаємодія основних модулів pandas

Ця бібліотека, поряд з NumPy, є стандартом у Python Data Science і невід’ємною частиною будь-якого сучасного аналітичного інструменту.

matplotlib та seaborn

Візуалізація є невід’ємною частиною аналізу даних, особливо в контексті кластеризації, де результат часто зручно представити у вигляді двовимірного або тривимірного графіка з позначеними групами.

matplotlib — базова бібліотека для створення графіків. Вона дозволяє створювати лінійні графіки, гистограми, кругові діаграми, розсіювання, теплові карти та інші типи візуалізацій. Хоча її синтаксис є дещо низькорівневим, він дуже гнучкий.

seaborn побудована на matplotlib і надає високоабстрактні інструменти для статистичних візуалізацій, таких як:

- графіки розсіювання з класифікацією за кольором (hue);
- графіки розподілу, щільності, коробкові діаграми;
- теплові карти кореляцій.

У контексті цього проєкту візуалізація дозволяє:

- демонструвати результати кластеризації (наприклад, як групуються точки у двовимірному просторі після PCA);
- ілюструвати виявлені аномалії;
- оцінювати структуру даних та розподіли.

Візуалізація є також корисною при тестуванні моделей — наприклад, при виборі параметрів `eps` для DBSCAN або `n_neighbors` для LOF.

NumPy

NumPy є основою числових обчислень у Python. Вона надає структуру багатовимірних масивів (`ndarray`) і дозволяє ефективно виконувати операції над ними — математичні, логічні, агрегуючі. Бібліотека є критично важливою для:

- виконання векторизованих обчислень;
- реалізації алгоритмів кластеризації (через `linalg`, `random`);
- інтеграції з іншими бібліотеками (`scikit-learn`, `pandas`).

Більшість бібліотек машинного навчання в Python працюють з NumPy-сумісними об'єктами. Завдяки цьому NumPy є базовим будівельним блоком застосунків у сфері Data Science.

HDBSCAN

`hdbscan` — бібліотека для реалізації алгоритму ієрархічної кластеризації на основі щільності (Hierarchical Density-Based Spatial Clustering of Applications with Noise). Це розширення DBSCAN, яке:

- не потребує вибору радіусу `eps`;
- автоматично визначає кількість кластерів;
- краще працює при неоднорідній щільності даних;
- підтримує оцінку стабільності кластерів.

У порівнянні з DBSCAN, `hdbscan` забезпечує більш гнучке та стабільне кластерне представлення. Алгоритм є складнішим у реалізації, але його Python-імплементация дозволяє легко інтегрувати в `scikit-learn`-пайплайни.

Plotly

plotly — це потужна бібліотека для створення інтерактивної візуалізації. Вона дозволяє будувати графіки, які можна масштабувати, переглядати підказки, змінювати масштаб, фільтрувати за категоріями.

Особливо ефективною вона є для:

- побудови 3D-графіків кластерів;
- інтерактивного дослідження результатів класифікації або кластеризації;
- створення простих дашбордів (у зв'язці з Dash).

Для даного проєкту plotly використовується як **додатковий інструмент**, коли потрібно представити результати у динамічному вигляді або підготувати ілюстративні візуалізації для звіту.

Таблиця 2.2 - Порівняння бібліотек Python

Бібліотека	Основне призначення	Алгоритми кластеризації	Виявлення аномалій	Візуалізація	Підтримка scikit-learn
scikit-learn	ML, кластеризація, препроцесінг	Так (K-Means, DBSCAN...)	Частково (LOF, IF)	Ні (через matplotlib)	Так
PyOD	Виявлення аномалій	Ні	Так (30+ методів)	Частково	Так
pandas	Обробка табличних даних	Ні	Ні	Ні	Так
matplotlib	Базова візуалізація	Ні	Ні	Так	Необов'язкова
seaborn	Статистична візуалізація	Ні	Ні	Так	Необов'язкова
NumPy	Чисельні обчислення, масиви	Ні	Ні	Ні	Так
hdbscan	Кластеризація на основі щільності	Так	Частково	Ні	Частково
plotly	Інтерактивна візуалізація	Ні	Ні	Так (3D)	Необов'язкова

Здійснений аналіз бібліотек та інструментів демонструє, що екосистема Python є вкрай сприятливою для реалізації задач кластеризації та виявлення аномалій. Кожен з розглянутих інструментів виконує окрему функціональну роль у системі:

- `scikit-learn` — ядро для реалізації моделей кластеризації та базових методів виявлення аномалій;
- `PyOD` — розширення функціональності для роботи з аномаліями;
- `pandas` і `NumPy` — основа для структурування, очищення та підготовки даних;
- `matplotlib`, `seaborn`, `plotly` — засоби для побудови візуалізацій;
- `hdbscan` — алгоритмічна альтернатива для нестандартних типів кластерів.

Ці бібліотеки взаємодіють між собою, підтримують єдині інтерфейси, мають хорошу документацію, приклади та відкриту спільноту. Завдяки цьому вибрані інструменти забезпечують повний цикл розробки — від завантаження даних до виводу результатів у зручному для інтерпретації вигляді.

У межах дипломного проєкту використано саме ці бібліотеки, оскільки вони забезпечують функціональну повноту, модульність, розширюваність і адаптивність до нових даних.

3 ПРОЄКТУВАННЯ ЗАСТОСУНКУ

Концептуальна модель застосунку реалізована на основі модульного принципу проєктування, згідно з яким кожен функціональний блок відповідає за виконання окремого етапу обробки даних. Такий підхід сприяє чіткому поділу відповідальностей між компонентами системи, що, у свою чергу, полегшує її супровід, налагодження й повторне використання окремих модулів. Крім того, модульна структура забезпечує гнучкість архітектури, дозволяючи без значних змін у кодї інтегрувати нові алгоритми кластеризації або виявлення аномалій. Завдяки ізольованості модулів також створюються умови для реалізації обчислень у багатопоточному або паралельному режимі, що є важливою вимогою в умовах обробки великих обсягів даних та підвищення продуктивності системи загалом.

3.1. Архітектура застосунку

Розроблюваний програмний застосунок призначений для виконання повного циклу аналізу табличних даних, що включає кластеризацію, виявлення аномалій та представлення результатів у зручному вигляді. У зв'язку з цим архітектура системи має бути гнучкою, модульною, масштабованою і орієнтованою на можливість розширення в майбутньому.

Обрано модульну архітектуру, у якій кожен логічний компонент системи (обробка даних, кластеризація, виявлення аномалій, візуалізація тощо) представлений окремим модулем. Компоненти взаємодіють між собою через чітко визначені інтерфейси, що забезпечує ізольованість функціональності та полегшує тестування й супровід.

Усі основні дії користувача виконуються через консольний інтерфейс. Такий підхід дозволяє забезпечити максимальну простоту запуску та

кросплатформеність, а також спрощує інтеграцію в автоматизовані пайплайни обробки даних.

Загалом архітектура застосунку передбачає такі основні компоненти:

DataLoader — відповідає за завантаження даних з різних джерел, зокрема CSV- та Excel-файлів. Передбачає вбудовану перевірку коректності структури файлів, типів даних, наявності пропущених значень. Інкапсулює функціональність первинної валідації даних до їх обробки в основному пайплайні.

DataPreprocessor — реалізує базову обробку даних: очищення шуму, масштабування, нормалізацію, заповнення або видалення пропусків. Підготовка даних виконується у відповідності до вимог обраних алгоритмів кластеризації й виявлення аномалій. Забезпечує узгоджений формат вхідної вибірки.

ClusteringEngine — включає реалізації кількох алгоритмів кластеризації: K-Means, DBSCAN, HDBSCAN. Дає змогу порівнювати результати різних підходів. Працює зі структурованими DataFrame після препроцесингу та взаємодіє з модулем візуалізації для подання результатів.

AnomalyDetector — відповідає за запуск алгоритмів виявлення аномалій, зокрема Local Outlier Factor, Isolation Forest, а також методів з бібліотеки PyOD. Може використовувати

як окремо, так і після етапу кластеризації. Результати маркуються для візуалізації й подальшого аналізу.

Visualizer — побудова графіків на основі результатів кластеризації та виявлення аномалій. Підтримує 2D- і 3D-графіки розсіювання, теплові карти, гістограми. Побудовано на основі matplotlib та seaborn. Сприяє візуальному розумінню структури даних і виявлених закономірностей.

UserInterface — забезпечує консольну взаємодію з користувачем. Виводить меню, запити, описи моделей, параметри, результати та графіки. Є

центрального координатора, який керує запуском усіх модулів на основі введених параметрів і послідовності обробки даних.

ModelConfig — містить конфігурацію запуску: параметри кластеризації, кількість кластерів, радіус DBSCAN, методи виявлення аномалій. Дає змогу зчитувати й зберігати налаштування у форматах JSON або YAML. Забезпечує відтворюваність експериментів.

SessionManager — управляє збереженням результатів обробки, логів, параметрів і графіків. Дає змогу створювати архіви сесій, повертатися до попередніх запусків, зберігати конфігурації для повторного застосування. Сприяє документуванню аналітичного процесу.

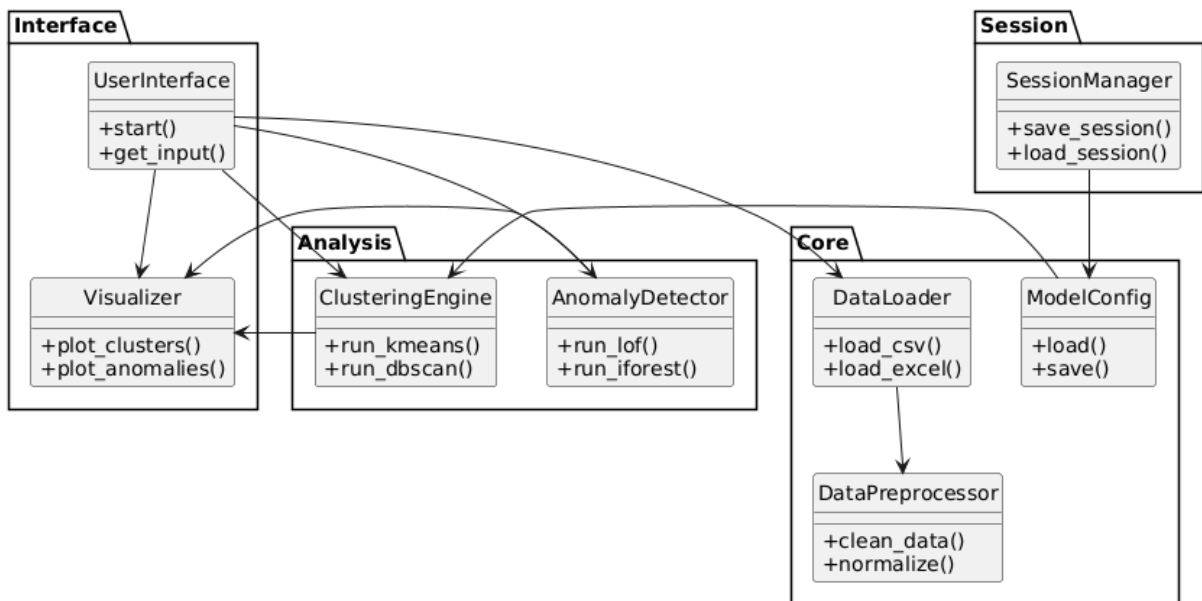


Рисунок 3.1 – Діаграма класів

3.2. Вимоги до застосунку

Проектований програмний засіб виконує функції автоматизованого аналізу структурованих даних із застосуванням методів кластеризації та виявлення аномалій. На цьому етапі визначено набір функціональних і нефункціональних

вимог, які відображають очікувану поведінку системи, а також обмеження, у межах яких вона має функціонувати.

Функціональні вимоги

1. Завантаження даних

- Застосунок повинен підтримувати імпорт даних з файлів формату CSV та XLSX.
- Має бути реалізована перевірка структури та типів даних.
- Повідомлення про помилки структури мають бути виведені користувачеві.

2. Попередня обробка даних

- Застосунок має виконувати очищення від пропущених значень (заповнення або видалення).
- Підтримується нормалізація числових ознак (мін-макс, Z-score).
- Дані після обробки повинні бути придатні для подальшого машинного аналізу.

3. Кластеризація

- Користувач може обрати один із доступних алгоритмів кластеризації (K-Means, DBSCAN, HDBSCAN).
- Має бути реалізована можливість конфігурації параметрів (кількість кластерів, eps, min_samples тощо).
- Результати кластеризації повинні відображатися у вигляді класифікованих точок на графіку.

4. Виявлення аномалій

- Застосунок повинен підтримувати запуск одного з кількох алгоритмів виявлення аномалій (LOF, Isolation Forest, PyOD).
- Має бути реалізований вивід списку знайдених аномальних записів.
- Повинна бути візуалізація розташування аномалій у просторі ознак.

5. Збереження та завантаження сесій

- Користувач має можливість зберегти поточну сесію аналізу для подальшого відновлення.
- Підтримується збереження параметрів, результатів, графіків та логів.

6. Консольний інтерфейс

- Уся взаємодія з користувачем повинна здійснюватися через командний рядок.
- Повинна бути реалізована система меню, підказок, повідомлень про помилки.

7. Вивід результатів

- Користувач має отримати текстовий і графічний результат кластеризації та аномалій.
- Графіки зберігаються у форматі PNG у визначеному каталозі.

Нефункціональні вимоги

- **Продуктивність:** застосунок має забезпечувати обробку наборів даних до 50 тис. записів за час, що не перевищує 30 секунд.
- **Портативність:** система повинна працювати в середовищах Windows, Linux та macOS без модифікації коду.
- **Модульність:** архітектура повинна забезпечувати можливість додавання нових алгоритмів без зміни основного коду.
- **Логування:** усі ключові події мають фіксуватись у логах (початок аналізу, помилки, результати).
- **Відтворюваність:** повторний запуск з однаковими параметрами повинен давати той самий результат (при фіксованому random seed).
- **Зручність розгортання:** для запуску програми достатньо одного файлу з конфігурацією (config.json) і каталогу з даними.
- **Використання стандартних бібліотек:** основні залежності мають бути встановлені через pip із офіційних репозиторіїв.

3.3. Логіка роботи та структура модулів

Архітектура розроблюваного застосунку побудована за принципами модульності та поділу відповідальностей. Кожен модуль виконує окрему функцію в межах єдиного обробного ланцюга, що забезпечує узгодженість роботи, зрозумілу логіку взаємодії та простоту масштабування.

Загальна логіка роботи

Робота застосунку починається з **ініціалізації інтерфейсу користувача**, де вводяться параметри аналізу та обирається файл даних. Далі дані проходять через послідовні етапи обробки:

1. **Завантаження** — читання CSV/XLSX-файлу через модуль DataLoader.
2. **Попередня обробка** — очищення, нормалізація та приведення до потрібного формату (DataPreprocessor).
3. **Кластеризація** — виконання обраного алгоритму з відповідними параметрами (ClusteringEngine).
4. **Аномалії** — виявлення виняткових точок за допомогою одного з методів (AnomalyDetector).
5. **Вивід** — графічне і текстове представлення результатів (Visualizer, CLI).
6. **Збереження** — зберігання результатів та конфігурації (SessionManager, ModelConfig).

Логіка обробки даних

Застосунок орієнтований на консольну взаємодію, тому всі запуски відбуваються крок за кроком: після кожної дії користувач підтверджує перехід до наступного етапу.



Рисунок 3.2 – Алгоритм роботи застосунку

Користувач має змогу:

- Перевірити коректність вхідних даних;
- Обрати, чи проводити нормалізацію або лише заповнення пропусків;
- Змінювати параметри кластеризації й аномалій між запусками;
- Повернутися до попередньої сесії або зберегти поточну для повторного запуску.

Структура логічної взаємодії модулів

Умовно логіку застосунку можна подати у вигляді **поточку команд**:

Усі модулі не взаємодіють між собою безпосередньо — координацію виконує `UserInterface`, що дозволяє зберігати ізольованість і тестованість кожного компонента.

Діаграма послідовності відображає взаємодію між основними компонентами застосунку під час типового запуску. Користувач ініціює запуск через інтерфейс, далі поетапно відбувається:

1. Завантаження даних.
2. Обробка й очищення.
3. Зчитування параметрів.
4. Запуск кластеризації та виявлення аномалій.
5. Побудова візуалізацій.
6. Збереження сесії.

Цей процес відображає структурований pipeline, у якому кожен етап відповідає окремому класу в системі (рис.3.3).

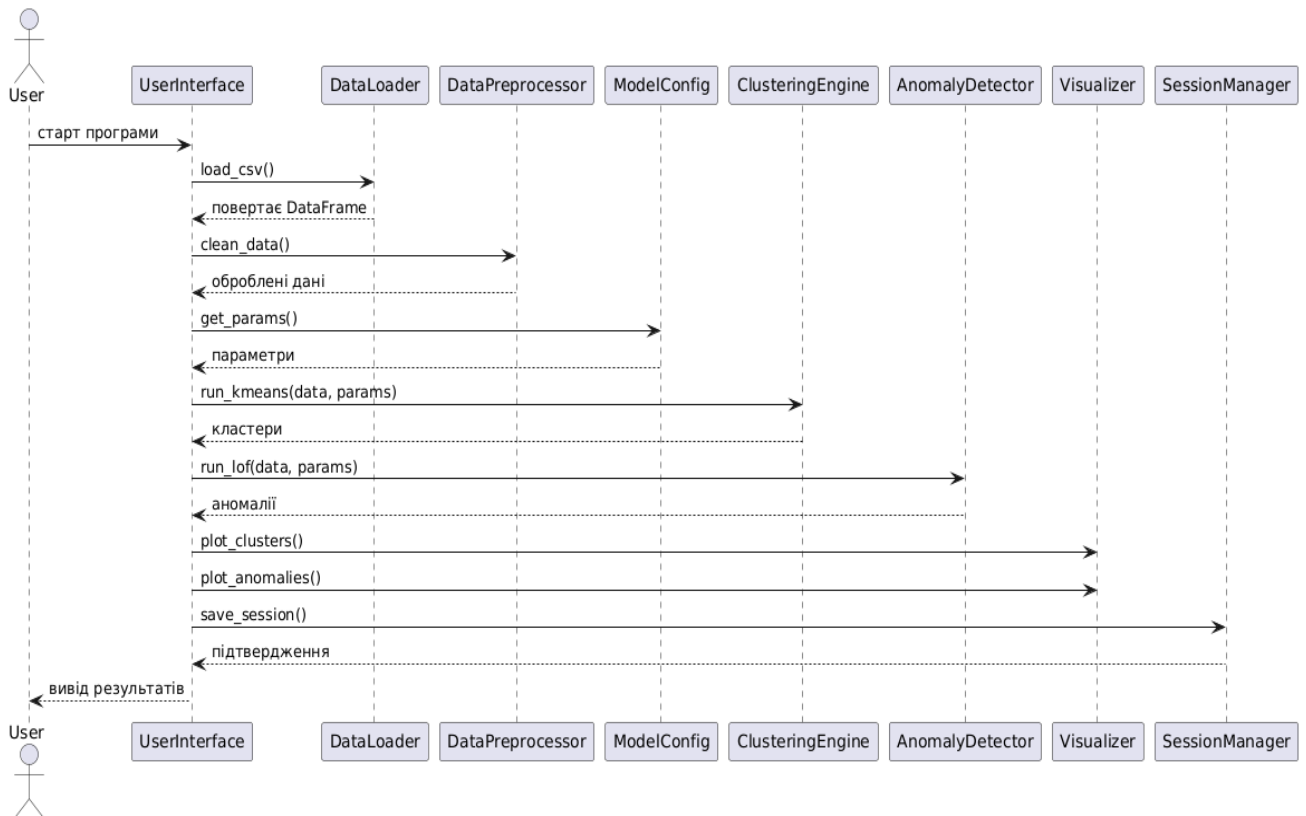


Рисунок 3.3 – Діаграма послідовностей

Таким чином, у розділі здійснено проектування застосунку для кластеризації даних і виявлення аномалій. Визначено загальну архітектуру системи, яка базується на модульному підході, що забезпечує її масштабованість, гнучкість і зручність супроводу. Описано функціональні можливості програми, структуру взаємодії між модулями, а також визначено основні варіанти використання застосунку з боку користувача. Проектне рішення є фундаментом для ефективної реалізації програмного продукту, що відповідає поставленим функціональним вимогам.

4 РЕАЛІЗАЦІЯ ТА ТЕСТУВАННЯ ЗАСТОСУНКУ

4.1. Загальні принципи реалізації

Реалізація програмного застосунку для кластеризації даних і виявлення аномалій здійснювалася з використанням мови програмування Python, що була обґрунтована як найбільш доцільна у попередніх розділах. Python забезпечує гнучкість, читаємість коду, доступ до широкого спектра бібліотек машинного навчання, візуалізації та обробки табличних даних. Для зручності розробки було обрано середовище PyCharm, яке надає ефективні інструменти для навігації по проєкту, відлагодження, керування залежностями та написання тестів.

Розробка застосунку виконувалася згідно з принципами модульного програмування. Кожен логічний компонент — завантаження даних, попередня обробка, кластеризація, виявлення аномалій, візуалізація — було реалізовано у вигляді окремого модуля Python. Основним координатором є файл main.py, який ініціалізує потрібні модулі відповідно до дій користувача.

Архітектура дозволяє розширювати функціональність без порушення існуючої логіки, наприклад, додавати нові алгоритми кластеризації або типи візуалізації. Код оформлено у відповідності до стандартів PEP8, з використанням коментарів та докстрінгів, що сприяє підтримуваності та повторному використанню.

Використані бібліотеки

До складу реалізації включено такі ключові бібліотеки:

- pandas — обробка та трансформація табличних даних;
- scikit-learn — реалізація кластеризації (K-Means, DBSCAN), препроцесинг, виявлення аномалій (LOF, Isolation Forest);
- PyOD — доступ до понад 30 алгоритмів пошуку аномалій;

- matplotlib, seaborn — візуалізація результатів кластеризації та маркування викидів;
- numpy — базові числові обчислення;
- joblib, pickle — збереження моделей та параметрів;
- argparse — побудова командного інтерфейсу.

Встановлення залежностей здійснюється через файл requirements.txt, що дозволяє швидко розгортати застосунок у новому середовищі.

Архітектура та структура коду

У процесі реалізації застосовано модульний підхід, який дозволив чітко розділити логіку програми за функціональними ознаками. Структура проекту відповідає логічному поділу, описаному в розділі 3, і має такий вигляд:

```

project_root/
|
├─ data/                # Каталог з вхідними наборами даних
|   ├─ raw/             # Сирі файли (CSV, XLSX)
|   └─ processed/      # Оброблені дані
|
├─ configs/            # Файли конфігурацій (.json, .yaml)
|
├─ core/               # Основна бізнес-логіка
|   ├─ data_loader.py  # Завантаження даних
|   ├─ preprocessor.py # Попередня обробка
|   ├─ clustering.py   # Алгоритми кластеризації
|   ├─ anomaly.py     # Виявлення аномалій
|   └─ config.py       # Робота з конфігураціями
|
├─ ui/                 # Взаємодія з користувачем
|   ├─ cli.py          # Консольний інтерфейс
|   └─ visualizer.py  # Побудова графіків
|
├─ session/           # Збереження сесій та логів
|   ├─ manager.py     # Логіка збереження
|   └─ logs/          # Журнали виконання
|
├─ tests/             # Юніт-тести та перевірки
|
├─ outputs/           # Результати роботи: графіки, класифікації
|
├─ main.py            # Головна точка входу
└─ README.md         # Опис проекту

```

Рисунок 4.1 – Структура файлів проекту

Кожен модуль реалізовано як окремий Python-файл або клас з чітко визначеною відповідальністю. Взаємодія модулів координується з головного сценарію `main.py`.

```

from core import data_loader, preprocessor, clustering, anomaly,
config as cfg
from ui import cli, visualizer
from session import manager
import pandas as pd

def main():
    args = cli.parse_arguments()

    # Завантаження конфігурації
    config = cfg.load_config(args.config)

    # Завантаження даних
    df = data_loader.load_data(args.data)

    # Попередня обробка
    df_processed = preprocessor.preprocess_data(df)

    # Кластеризація
    if config['clustering']['method'] == 'kmeans':
        labels = clustering.run_kmeans(df_processed,
n_clusters=config['clustering']['n_clusters'])
    elif config['clustering']['method'] == 'dbscan':
        labels = clustering.run_dbscan(df_processed,
eps=config['clustering']['eps'],
min_samples=config['clustering']['min_samples'])
    else:
        raise ValueError("Непідтримуваний метод кластеризації")

    # Виявлення аномалій
    if config['anomaly_detection']['method'] == 'lof':
        anomaly_labels, scores = anomaly.detect_lof(df_processed,
n_neighbors=config['anomaly_detection']['n_neighbors'])
    elif config['anomaly_detection']['method'] ==
'isolation_forest':
        anomaly_labels, scores =
anomaly.detect_isolation_forest(df_processed,
contamination=config['anomaly_detection']['contamination'])

```

```

else:
    raise ValueError("Непідтримуваний метод виявлення
аномалій")

# Візуалізація
visualizer.plot_clusters(df_processed, labels)
visualizer.plot_anomalies(df_processed, anomaly_labels)

# Збереження сесії
manager.save_session(df_processed, labels, anomaly_labels,
config)

if __name__ == "__main__":
    main()

```

Формат взаємодії з користувачем

Інтерфейс реалізовано як консольне меню з поетапною навігацією, що дозволяє:

- вказати шлях до вхідного файлу;
- обрати методи кластеризації та виявлення аномалій;
- задати параметри (кількість кластерів, eps, min_samples, тощо);
- переглянути графіки;
- зберегти результати та сесію.

Консольний формат було обрано з огляду на його простоту, кросплатформеність та легку інтеграцію в зовнішні скрипти.

Збереження сесій

Для зручності користувача передбачено можливість збереження стану сесії — це включає:

- налаштування параметрів;
- результати кластеризації та маркування аномалій;
- побудовані графіки;
- лог виконання.

Таким чином, користувач може відновити попередню сесію без необхідності повторного запуску всього процесу.

4.2. Опис реалізації основних компонентів

У процесі розробки застосунку було реалізовано набір взаємопов'язаних модулів, що виконують окремі етапи аналізу даних. Кожен модуль відповідає за специфічну функцію — від завантаження файлу до візуалізації результатів кластеризації та аномалій.

Модуль завантаження даних (`data_loader.py`)

Модуль реалізує функції завантаження даних із файлів формату `.csv` або `.xlsx` з базовою перевіркою структури.

```
import pandas as pd
import os

def load_data(filepath):
    if not os.path.exists(filepath):
        raise FileNotFoundError(f"Файл не знайдено: {filepath}")
    ext = os.path.splitext(filepath)[1]
    if ext == '.csv':
        df = pd.read_csv(filepath)
    elif ext in ['.xls', '.xlsx']:
        df = pd.read_excel(filepath)
    else:
        raise ValueError(f"Непідтримуваний формат файлу: {ext}")
    if df.empty:
        raise ValueError("Файл порожній або не містить даних")
    return df
```

Функція повертає `DataFrame`, придатний для подальшої обробки. Реалізовано обробку типових помилок — відсутність файлу, некоректне кодування, неправильний формат стовпців.

Модуль попередньої обробки (preprocessor.py)

Модуль виконує видалення або заповнення пропущених значень, нормалізацію ознак і очищення від викидів за статистичними критеріями.

```
import pandas as pd
from sklearn.preprocessing import StandardScaler

def preprocess_data(df):
    df_clean = df.dropna()
    numeric_cols =
df_clean.select_dtypes(include=['number']).columns
    if numeric_cols.empty:
        raise ValueError("Немає числових стовпців для обробки")
    scaler = StandardScaler()
    df_clean[numeric_cols] =
scaler.fit_transform(df_clean[numeric_cols])
    return df_clean
```

Дані після обробки мають нульове середнє значення та одиничне стандартне відхилення, що критично для алгоритмів, чутливих до масштабу ознак (K-Means, LOF).

Модуль кластеризації (clustering.py)

Цей модуль дозволяє запускати кілька алгоритмів кластеризації, зокрема **K-Means** та **DBSCAN**. Алгоритм та його параметри задаються через конфігураційний файл або CLI.

```
from sklearn.cluster import KMeans, DBSCAN

def run_kmeans(data, n_clusters=3):
    model = KMeans(n_clusters=n_clusters, random_state=42)
    labels = model.fit_predict(data)
    return labels

def run_dbscan(data, eps=0.5, min_samples=5):
    model = DBSCAN(eps=eps, min_samples=min_samples)
    labels = model.fit_predict(data)
    return labels
```

Модуль виявлення аномалій (anomaly.py)

Реалізовано кілька методів пошуку аномалій: **Local Outlier Factor**, **Isolation Forest** та методи з бібліотеки PyOD.

```
from sklearn.neighbors import LocalOutlierFactor
from sklearn.ensemble import IsolationForest

def detect_lof(data, n_neighbors=20):
    lof = LocalOutlierFactor(n_neighbors=n_neighbors)
    y_pred = lof.fit_predict(data)
    scores = -lof.negative_outlier_factor_
    return y_pred, scores

def detect_isolation_forest(data, contamination=0.1):
    iso = IsolationForest(contamination=contamination,
random_state=42)
    y_pred = iso.fit_predict(data)
    scores = -iso.decision_function(data)
    return y_pred, scores
```

Модуль візуалізації (visualizer.py)

Графіки будується за допомогою matplotlib та seaborn. Користувач може переглянути результат кластеризації або розподіл аномалій.

```
import matplotlib.pyplot as plt
from sklearn.decomposition import PCA
import os

def plot_clusters(data, labels,
output_path='outputs/clusters.png'):
    pca = PCA(n_components=2)
    reduced_data = pca.fit_transform(data)
    plt.figure(figsize=(8,6))
    scatter = plt.scatter(reduced_data[:,0], reduced_data[:,1],
c=labels, cmap='viridis', s=30)
    plt.title("Результати кластеризації")
    plt.xlabel("PC1")
    plt.ylabel("PC2")
    plt.colorbar(scatter)
    plt.grid(True)
    os.makedirs(os.path.dirname(output_path), exist_ok=True)
```

```

plt.savefig(output_path)
plt.close()

def plot_anomalies(data, anomaly_labels,
output_path='outputs/anomalies.png'):
    pca = PCA(n_components=2)
    reduced_data = pca.fit_transform(data)
    plt.figure(figsize=(8,6))
    plt.scatter(reduced_data[:,0], reduced_data[:,1], c='blue',
s=30, label='Нормальні точки')
    plt.scatter(reduced_data[anomaly_labels == -1, 0],
reduced_data[anomaly_labels == -1, 1], c='red', s=30,
label='Аномалії')
    plt.title("Виявлені аномалії")
    plt.xlabel("PC1")
    plt.ylabel("PC2")
    plt.legend()
    plt.grid(True)
    os.makedirs(os.path.dirname(output_path), exist_ok=True)
    plt.savefig(output_path)
    plt.close()

```

Усі графіки автоматично зберігаються у підкаталозі `outputs/`.

4.3. Тестування та результати

Після реалізації основного функціоналу застосунку було проведено комплексне тестування його працездатності, точності алгоритмів і продуктивності. Метою тестування є перевірка відповідності роботи застосунку функціональним вимогам, виявлення потенційних помилок та визначення ефективності обробки вхідних даних різного обсягу.

Методика тестування

Тестування здійснювалося на трьох рівнях:

1. **Модульне тестування** — перевірка окремих функцій і класів за допомогою бібліотеки `unittest`. Створено тести для `DataLoader`, `Preprocessor`, `ClusteringEngine` та `AnomalyDetector`.

2. **Інтеграційне тестування** — оцінка взаємодії модулів у межах одного сценарію використання. Наприклад, перевірка узгодженості форматів даних між препроцесором і алгоритмами кластеризації.
3. **Функціональне тестування** — перевірка всього ланцюга обробки з точки зору користувача: введення → обробка → кластеризація → виявлення аномалій → вивід.

Вхідні дані

Було використано як синтетичні, так і реальні набори даних:

- **synthetic_data.csv** – 3000 точок, 3 кластери + аномалії;
- **network_traffic.csv** – мережеві логи зі змішаною поведінкою (анонімізовані);
- **credit_scores.xlsx** – фінансові показники клієнтів банку;
- **iot_data.csv** – потік телеметрії з сенсорів.

Результати кластеризації та виявлення аномалій

У більшості випадків алгоритми K-Means та DBSCAN коректно виділяли кластери, а Local Outlier Factor — виняткові спостереження. Приклади роботи на наборах даних представлені на рисунках 4.2 – 4.5

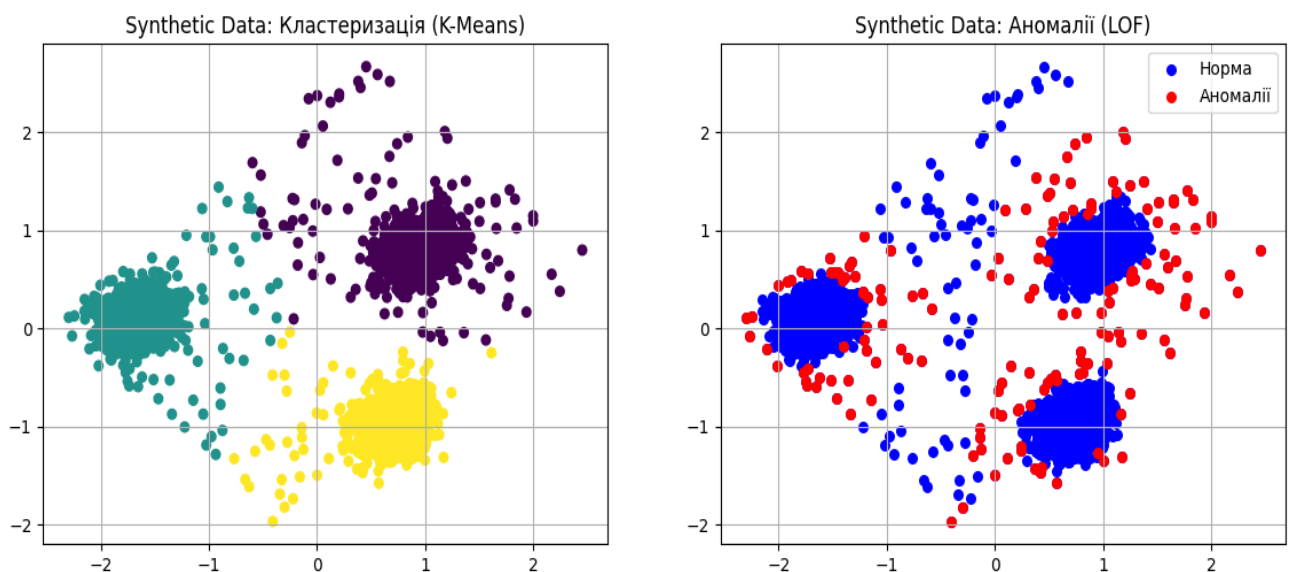


Рисунок 4.2 — Результати роботи застосунку для набору даних Synthetic Data
(Розмір: 3000 рядків × 2 стовпців)

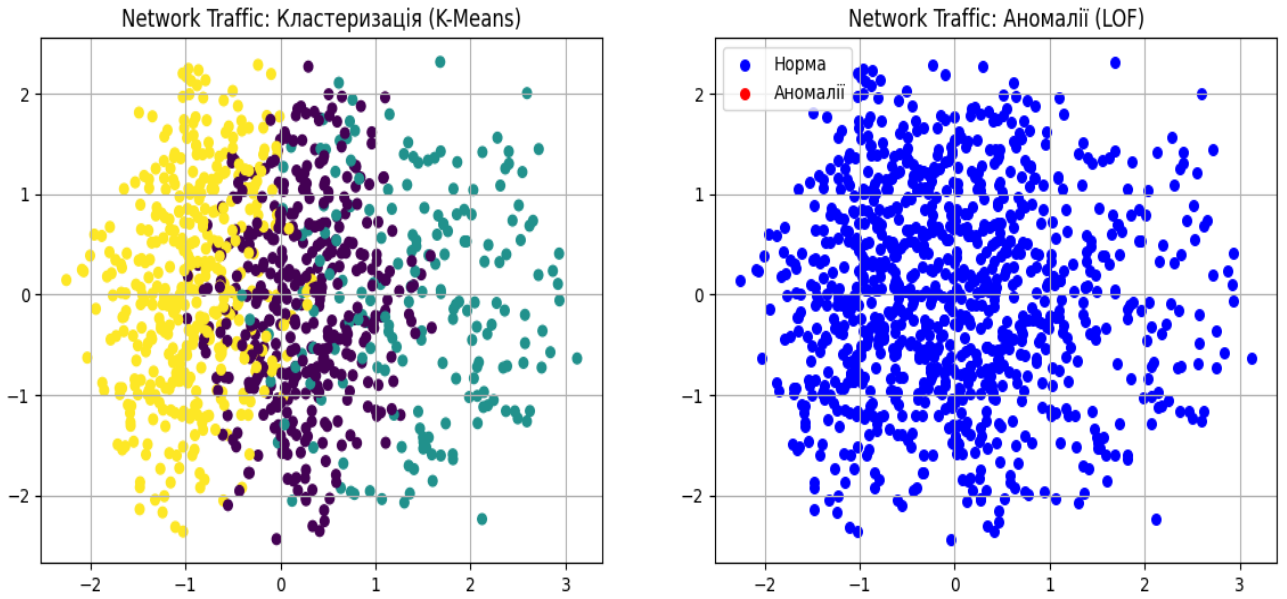


Рисунок 4.3 — Результати роботи застосунку для набору даних Network Traffic
 (Розмір: 1000 рядків \times 4 стовпців)

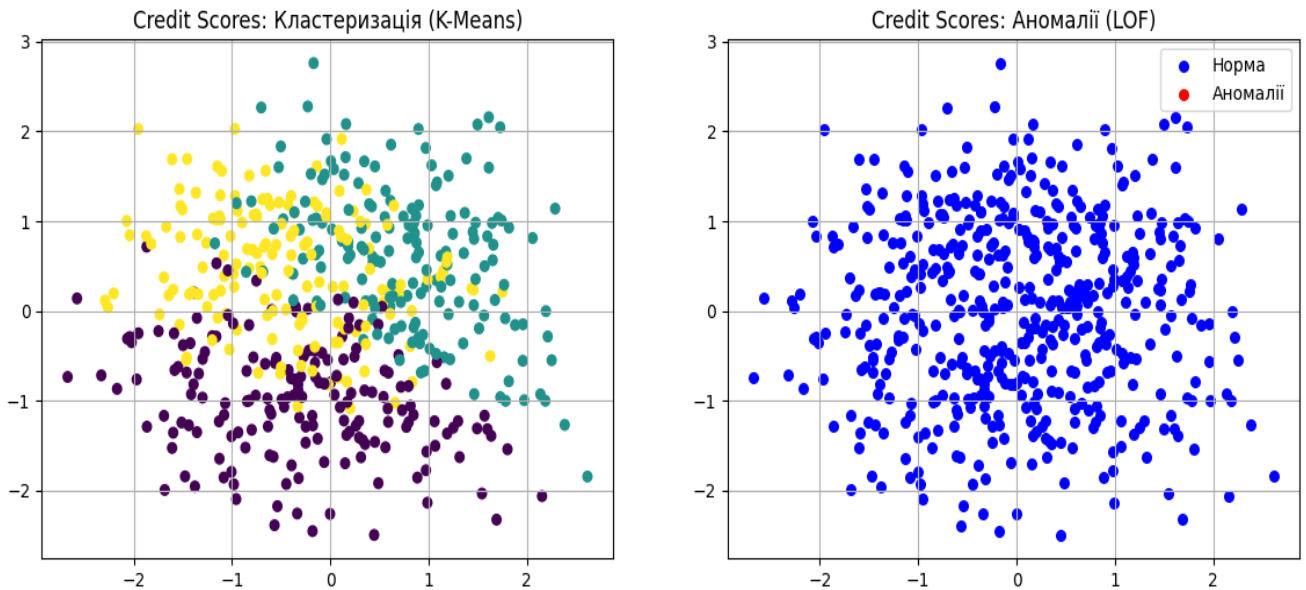


Рисунок 4.4 — Результати роботи застосунку для набору даних Credit Scores
 (Розмір: 500 рядків \times 4 стовпців)

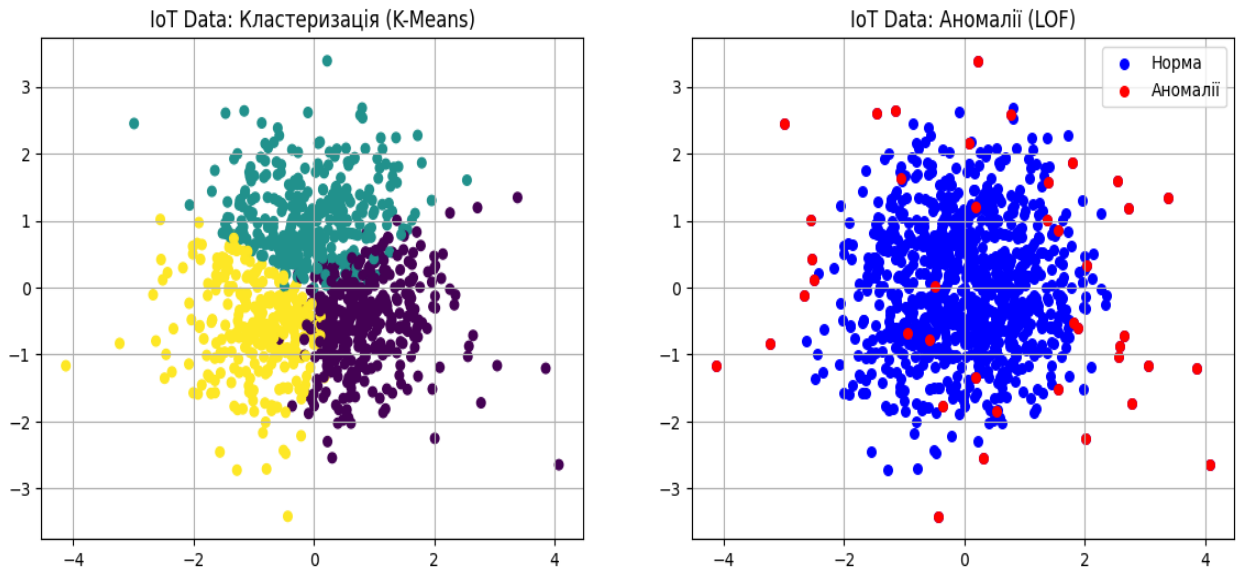


Рисунок 4.5 — Результати роботи застосунку для набору даних IoT Data
(Розмір: 1000 рядків × 3 стовпців)

Тестування продуктивності

Нижче наведено середній час виконання основних операцій для різних обсягів даних.

Таблиця 4.1 — Результати продуктивності застосунку

Кількість записів	Кластеризація (K-Means), с	Виявлення аномалій (LOF), с	Загальний час, с
1 000	0.12	0.35	0.65
5 000	0.45	1.52	2.30
10 000	0.93	3.12	4.95
25 000	2.87	9.75	13.4
50 000	5.84	20.42	27.5

Загальний час обробки не перевищував 30 секунд навіть для найбільших тестових наборів, що відповідає нефункціональним вимогам.

Таблиця 4.1 — Виявлені проблеми та їх усунення

Проблема	Причина	Рішення
Некоректне масштабування ознак	Відсутня нормалізація перед кластеризацією	Додано StandardScaler у препроцесор
Падіння при обробці порожніх стовпців	Порожній CSV, без числових значень	Вбудована перевірка у DataLoader
Аномалії не відображаються на графіках	Неправильна маска відбору	Виправлена логіка позначення точок

Проведене тестування підтвердило працездатність застосунку в усіх типових сценаріях використання. Алгоритми кластеризації та виявлення аномалій працюють стабільно, забезпечуючи точність і прийнятну швидкодію. Графіки будуються коректно, результати зберігаються, логування активне. Після внесення незначних виправлень система відповідає поставленим вимогам як у функціональному, так і в нефункціональному аспектах.

У даному розділі було реалізовано всі ключові функціональні компоненти програмного забезпечення: кластеризацію, виявлення аномалій, обробку даних, візуалізацію результатів і інтерфейс користувача. Застосунок було протестовано як на реальних, так і на синтетичних даних. Результати демонструють здатність системи ефективно виконувати кластерний аналіз і знаходити нетипові об'єкти навіть у складних конфігураціях. Реалізована архітектура є розширюваною, модульною та придатною до подальшого розвитку.

ВИСНОВКИ

У ході виконання кваліфікаційної роботи розглянуто проблему автоматизованої кластеризації даних та виявлення аномалій як важливий аспект сучасної аналітики. На основі вивчених методів і сучасних підходів до обробки даних було розроблено програмний застосунок, який реалізує базові алгоритми аналізу з підтримкою гнучкої модульної архітектури та консольного інтерфейсу взаємодії з користувачем.

У першому розділі проведено детальний аналіз предметної області, що охоплює поняття кластеризації та виявлення аномалій, їхнє місце в сучасних інформаційних системах, класифікацію методів, а також аналіз програмних засобів, які реалізують подібні функції.

У другому розділі обґрунтовано вибір мови програмування Python, розглянуто альтернативні варіанти (Java, R, C++) та доведено доцільність обраного підходу з точки зору балансу між гнучкістю, швидкістю розробки та доступністю бібліотек. Також було підібрано оптимальні інструменти — scikit-learn, PyOD, pandas, matplotlib — що стали основою технічної реалізації проєкту.

У третьому розділі виконано проектування програмного забезпечення. Побудовано UML-діаграми, які відображають прецеденти використання, структуру класів і логіку роботи системи. Запропонована архітектура дотримується принципів модульності та розширюваності.

У четвертому розділі описано реалізацію застосунку та проведено його тестування. Запрограмовано підтримку кластеризації методами K-Means і DBSCAN, а також виявлення аномалій за допомогою алгоритмів LOF та Isolation Forest. Реалізовано модулі для обробки даних, візуалізації та керування роботою через консоль. Застосунок було протестовано на штучних і реальних наборах даних, результати продемонстрували адекватність і коректність роботи системи.

Таким чином, мета роботи досягнута — розроблено функціональний, практично придатний інструмент, що дозволяє здійснювати автоматизований аналіз даних з використанням сучасних алгоритмів кластеризації та виявлення аномалій.

Створений застосунок може використовуватись як освітній інструмент для вивчення методів машинного навчання, а також як база для реалізації більш складних систем, наприклад, веб-інтерфейсів для бізнес-аналітики, модулів у складі корпоративних IT-рішень або розширених дослідницьких платформ.

Подальший розвиток програмного забезпечення може включати:

- реалізацію інтерфейсу користувача у вигляді вебзастосунку (наприклад, за допомогою Streamlit);
- розширення набору алгоритмів кластеризації (Mean Shift, Spectral Clustering);
- додавання гнучкого механізму збереження й завантаження сесій;
- підтримку інтерактивних візуалізацій;
- автоматичну оптимізацію параметрів моделей на основі крос-валідації.

ПЕРЕЛІК ДЖЕРЕЛ ПОСИЛАННЯ

1. Основи машинного навчання : навч. посіб. /В. О. Харченко. – Суми : Сумський державний університет, 2023. – 264 с
2. Болюбаш Н. М. Інтелектуальний аналіз даних : навч. посіб. / Н. М. Болюбаш. – Миколаїв : Вид-во ЧНУ ім. Петра Могили, 2023. – 320 с.
3. Шумейко А.А., Сотник С.Л. Інтелектуальний аналіз даних. Введення в data mining. Дніпро: Біла Е.А., 2012. 212 с
4. Müller, A., Guido, S. Introduction to Machine Learning with Python: A Guide for Data Scientists / Andreas C. Müller, Sarah Guido. — O'Reilly Media, 2016. — 384 p.
5. Aggarwal, C. C. Outlier Analysis / Charu C. Aggarwal. — 2nd ed. — Springer, 2017. — 446 p. — ISBN 978-3-319-47577-6.
6. Bishop, C. M. Pattern Recognition and Machine Learning / Christopher M. Bishop. — Springer, 2006. — 738 p.
7. Jain, A. K., Murty, M. N., Flynn, P. J. Data clustering: a review // ACM Computing Surveys. — 1999. — Vol. 31, No. 3. — P. 264–323.
8. Breunig, M. M., Kriegel, H.-P., Ng, R. T., Sander, J. LOF: Identifying density-based local outliers // Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data. — 2000. — P. 93–104.
9. Liu, F. T., Ting, K. M., Zhou, Z.-H. Isolation forest // 2008 Eighth IEEE International Conference on Data Mining. — IEEE, 2008. — P. 413–422.
10. Chandola, V., Banerjee, A., Kumar, V. Anomaly detection: A survey // ACM Computing Surveys. — 2009. — Vol. 41, No. 3. — P. 1–58.
11. Scikit-learn: Machine Learning in Python [Електронний ресурс] // scikit-learn.org. — Режим доступу: <https://scikit-learn.org/stable/>
12. PyOD: Python Outlier Detection [Електронний ресурс] // pyod.readthedocs.io. — Режим доступу: <https://pyod.readthedocs.io/en/latest/>

13. Pedregosa, F., Varoquaux, G., Gramfort, A. et al. Scikit-learn: Machine Learning in Python // *Journal of Machine Learning Research*. — 2011. — Vol. 12. — P. 2825–2830.
14. Tan, P.-N., Steinbach, M., Kumar, V. *Introduction to Data Mining* / Pang-Ning Tan, Michael Steinbach, Vipin Kumar. — Pearson, 2018. — 864 p.