

РОЗРОБКА DATA MESH ПЛАТФОРМИ ДАНИХ З ML ДОМЕНОМ АНАЛІЗУ ДАНИХ

М. Фостяк, Л. Демків

*Львівський національний університет імені Івана Франка,
Драгоманова, 50, 79005 Львів, Україна*

lidia.demkiv@gmail.com

У роботі створено модель data mesh сітки даних, яка забезпечує децентралізовану та розподілену архітектуру даних для проекту з надання фінансових послуг клієнтам. Запропоновано схему моделі з потоковими доменами на основі даних отриманих з Open Bank API. Реалізовано схему взаємодії та розподілу даних між доменами для створення маркетингових дата-продуктів, а також дата-продуктів ефективності компанії, моніторингу користувачів та інших.

Рівень аналізу даних включає домен ML моделей. У домені реалізовано методи класифікації даних. Встановлено класифікатори та параметри, за яких досягається найвища точність класифікації з врахуванням аугментації даних. Визначено параметри важливості даних та поведено аналіз результатів класифікації. Запропоновано враховувати результати ML аналізу даних для підвищення точності класифікації клієнтів, аналізу фінансових кредитних ризиків та визначення оптимальної відсоткової ставки.

Ключові слова: моделі зберігання даних, data mesh, домени даних, класифікація, ML аналіз даних.

Вступ

Проектування розвитку бізнесу на основі аналізу даних є актуальним напрямом, який проявляється у побудові різноманітних за функціональністю моделей опрацювання даних. Тренди сьогодення констатують, що навіть невеликі компанії можуть мати значні обсяги даних, особливо з урахуванням зростання цифрового бізнесу та доступу до різних джерел даних. Розмір компаній не завжди прямопропорційно корелює з обсягом даних, які вони обробляють. Малі і великі бізнеси можуть накопичувати значну кількість інформації з різних джерел, що приводить до необхідності забезпечення структурованого зберігання даних з метою їх наступного використання для розвитку бізнесу. Створення концептуального представлення даних та забезпечення зв'язків між даними всередині організації вимагає побудови певних моделей, які забезпечать своєчасне прийняття рішень на основі аналізу даних, що зберігаються у моделях, та мають вирішальне значення для збереження конкурентоспроможності бізнесу. Розвиток моделей сховищ даних від data warehouse (DWH) до data lake і до концепції data mesh представляє собою еволюцію підходів до зберігання, обробки і аналізу даних. Загалом, ці три підходи до управління даними представляють різні етапи в еволюції від централізованих та структурованих до розподілених та гнучких моделей. Кожен з них може бути корисним залежно від потреб компанії та особливостей даних.

Data warehouse (DWH) є централізованим сховищем даних, у якому дані зберігаються в структурованому форматі, найчастіше за схемою зірочки. Цей підхід передбачає одну систему зберігання та загальний доступ до даних. Звичайно, дані в DWH використовуються для звітності, аналітики та бізнес-інтелекту, але обробка даних, зазвичай, є централізованою і питання до даних вже передбачені наперед [1]. Будь яка необхідність змінити питання приводить до необхідності перебудови сховища. Аналітики даних не достатньо швидко можуть сформулювати відповідь на нові аналітичні питання керівництва та власників продуктів через необхідність модифікації запитів або перебудови процесів збереження даних. Тому наступною моделлю сховищ даних став data lake: Data lake є розподіленим сховищем даних, де дані зберігаються у сирому, неструктурованому або напівструктурованому форматі. Data lake спроектований для масштабованості та гнучкості, і дозволяє аналізувати дані у різних контекстах та використовувати різні інструменти для обробки даних. Однак, аналіз цих даних може бути складним через відсутність заздалегідь встановлених схем і структур. Це може ускладнити пошук, фільтрацію та вилучення потрібної інформації, а також вимагатиме значних зусиль для налагодження безпеки, моніторингу та підтримки інфраструктури.

Децентралізовані сховища data mesh

Data mesh - це підхід, який прагне забезпечити децентралізацію та самоорганізацію в управлінні даними [2]. Замість централізованого підходу до управління даними, Data mesh визначає області відповідальності для окремих команд або "доменів", які відповідають за власні дані та їх управління. Кожен домен може мати свої власні сховища даних (може бути як DWH, так і data lake), а також власні процеси обробки даних та інструменти аналізу. Основне завдання сітки даних - це сформулювати відповідні дані для конкретного підрозділу. Таке завдання можна вирішити шляхом створення доменів даних для відповідних даних продуктів та налагодження оптимальної взаємодії між доменами в процесі постачання даних продукту до команди споживача. Для формування сітки даних необхідно оцінити кількість необхідних доменів даних, а також кількість даних продуктів та функціональних зв'язків між ними, які необхідно реалізувати для команд (маркетинг, клієнти, продажі), що керують прийняттям рішень на основі даних.

Базовими типами доменів є операційні, аналітичні, інтеграційні. Операційні домени забезпечують доступ до операційних даних. Аналітичні домени включають дані, які використовуються для аналізу та прийняття стратегічних рішень. Інтеграційні відповідають за інтеграцію даних між різними доменами та зовнішніми системами і містять дані про API, інтеграційні точки, трансформації даних, протоколи обміну даними. Реальні data mesh структури розробляють власні домени для конкретних виробничих потреб. Команди доменів визначають, які дані необхідні для вирішення конкретної задачі або для надання певної аналітичної інформації, вивчають доступні джерела даних та визначають, які дані можна використовувати для створення продукту, отримують оперативні дані та проводять очищення, трансформацію, стандартизацію та інтеграцію даних з різних джерел, а також створюють аналітичні моделі даних, як продукти даних, для виконання власного аналізу. Таким чином операційні дані використовуються для створення даних продуктів, які сприяють раціональному прийняттю рішень або розумінню бізнес-процесів. Команда домену може публікувати та передавати продукти даних для задоволення потреб інших доменів у даних. Однак при

такому підході потрібно розуміти, що одні і ті ж дані дублюються для різних доменів. Цей недолік компенсується повною відповідальністю команди домену за дата продукт.

Доменна data mesh модель сховища даних фінансової структури

На рис. 1 зображено схема data mesh з трьома вхідними доменами А, Б, В. Кожен з цих доменів має своє операційне джерело даних. Операційні джерела даних можуть використовувати різні технології для зберігання даних та їх опрацювання, наприклад, одна частина може бути реляційною, а інша NoSQL базою даних. Далі кожна доменна команда будує дата продукт для свого домену (Дата-продукт Домену А, Б, В), в який включає лише очищені, опрацьовані та вибрані дані. Побудова дата продуктів передбачає слідування структурованому процесу створення моделей або конвеєрів даних, забезпечення автоматизованості та масштабованості даних, забезпечення механізму моніторингу даних, перевірки на відсутність або неправильність даних, обмеження доступу до конфіденційних даних, формування документації та керування життєвим циклом обробки даних. Таким чином у результаті отримуємо спрямовані на джерело дата продукти. Ці доменні дата продукти потім об'єднуються в комплексний агрегат, який містить всеохопні дані про сутності в системі. Далі створюємо спрямовані на користувача дата-продукти. Маркетинговий дата-продукт використовує продукт Домену А, а також зовнішні дані аналітики, з таких систем як Google Analytics чи Hotjar. Дата-продукт ефективності компанії, який використовує керівництво компанії для загальних метрик успішності, побудовано на агрегаті системи, адже споживач потребує якнайширшої картини стану справ.

Дата продукт для тренування моделей машинного навчання використовує дані з Доменів Б та В і надає трансформовані дані. Таким чином аналітики даних можуть зосередитися на розробці, навчанні та вдосконаленні своїх моделей, а не на підготовці даних. Власне на даних такого дата продукту і відбувалося дослідження алгоритмів далі в статті.

Параметри даних для ML моделей

Для побудови автоматизованих інтелектуальних систем розподілу ймовірностей повернення кредитів клієнтами на основі ML моделей необхідно зібрати та сформувати таблиці даних, які отримані з банківських реквізитів, а також даних про клієнтів та багато іншої інформації, пов'язаної з клієнтами та кредитуванням. У роботі використано модель даних, яка узгоджена з параметрами даних, що пропонується сервісом Open Bank API. Open Bank API – це відкритий банк, який забезпечує можливість побудови інноваційних додатків і послуг, адаптованих до фінансових даних користувачів [3]. Інформація про клієнта, яку можна отримати від Open Bank API, може включати: особисті дані, які клієнт надав під час реєстрації або в процесі взаємодії з банком, дані про баланс рахунків, історію транзакцій, категорії витрат, платежі, перекази коштів, статистику витрат та доходів, інформацію про поточні кредити, позики, графік виплат, інформацію про статус платежів (наприклад, чи були вони виконані, чи є вони в очікуванні тощо) та історію платежів, аналітичних звітів та статистичних даних щодо фінансових транзакцій та витрат клієнта. Формування даних для побудови ML моделі кредитного рейтингу проводилось після отримання реальних даних про щоденні детальні транзакції клієнтів, які відповідають інформації, яку клієнт отримує на своїй банківській картці. Дані були категоризовані, агреговані та анонімізовані. В таблиці 1 представлено назви отриманих колонок даних за трьома розділами: категоризовані

витрати, ризиковані витрати та категоризовані надходження. Характерною особливістю даних для кредитного рейтингу є їхня незбалансованість за колонкою Fail. Існують різні способи для перетворення незбалансованих даних у збалансовані [4]. В роботі використано підхід – аугментація даних [5]. Цей підхід також дасть відповідь на питання про те, чи буде отримано кращу точність ML алгоритмів для датасету, який тренується на більшій кількості згенерованих даних і тестується на реальних даних.

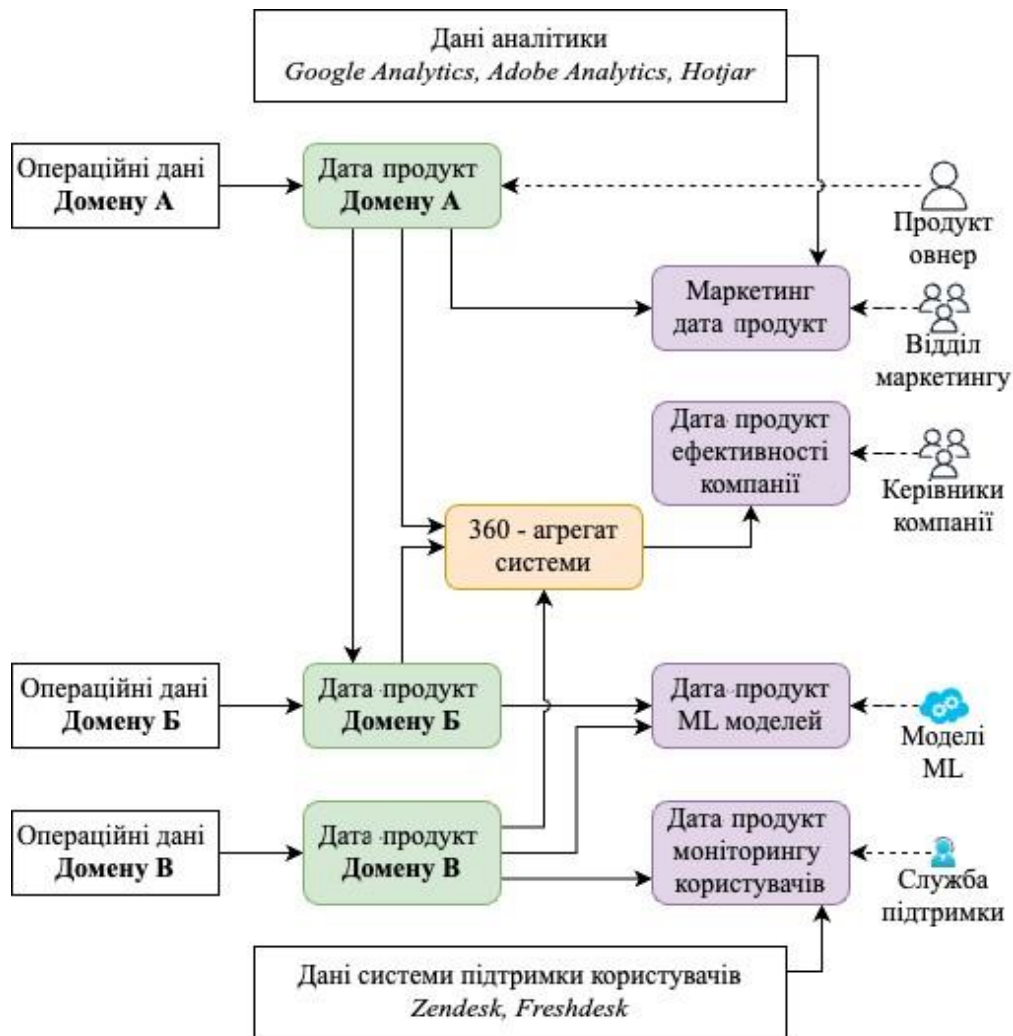


Рис.1. Схема структури data mesh
Fig. 1. Scheme of the data mesh structure

Для генерації даних використано Generative Adversarial Network (GAN), яка складається з двох частин: генератора та дискримінатора, а також ydata-synthetic —

бібліотеки з відкритим кодом для створення синтетичних даних. Генератор та дискримінатор – це дві нейромережі, які конкурують одна з одною за принципом: генератор створює зміст, який складно відрізнити від реального, а дискримінатор намагається відрізнити створений зміст від справжнього. В результаті отримано три набори даних: реальний незбалансований (дані 1), згенерований незбалансований (дані 2) та згенерований збалансований (дані 3).

Таблиця 1. Категоризований опис даних
Table 1. Description of categorized data

	Категоризовані витрати	Ризиковані витрати	Категоризовані надходження
	Кредитні витрати <i>Скільки коштів в місяць витрачає клієнт на оплату існуючих кредитів</i> -0,02 CreditPayments	Витрати на оплай казино та букмекерські контори 0,39 Gambling	Надходження -0,45 TotalInflow
	Витрати на продукти 0 Groceries	Неоплачені регулярні платежі <i>Комунальні платежі, платежі за моб зв'язок</i> 0,45 BouncedDirectDebits	Регулярні надходження категоризовані як зарплата -0,52 Income
	Витрати на транспорт <i>Громадський транспорт, таксі</i> -0,02 Commute	Кількість нарахувань за рахунок кредитних коштів -0,07 CreditMoney Inflows	Середньовідсоткове зростання зарплати за останній рік 0,0 SalaryGrowth
	Витрати на авто <i>Заправки, автосервіси, магазини автотоварів</i> -0,04 Car	Кількість коштів знятих колекторськими агенціям 0,4 CreditCollection Agencies	Надходження в рахунок інвестицій та збережень <i>Депозити, виплати дивідендів</i> 0,07 InvestmentIncome
	Витрати на подорожі <i>Квитки на поїзд, літак, паром. Готелі.</i> 0,01 Travel,	Дані отримані шляхом обчислення різниці між надходженнями і витратами	
	Зняття готівки -0,07 CashWithdrawal	Різниця надходження-витрати -0,52 InflowOutflowsDifference	
		Карткові перекази -0,02 CardTransfersOut -0,03 CardTransfersIn	
	BNPL платежі -0,01 Виплати за розтернінуванням	Дані про клієнта з Домену В 0,37 EmploymentIndustry 1 Fail	

Проведено дослідницький аналіз даних, у тому числі кореляційний, статистичний, візуальний, а також стандартизацію даних. Для даних 2 числові значення відповідних

коефіцієнтів кореляції отримано з матриці кореляції між колонками значень і Fail. Ці значення проказані у відповідних комірках Таблиці 1. Як видно з таблиці 1 найбільші позитивні значення кореляції з Fail отримано для параметрів BouncedDirectDebits, CreditCollectionAgencies, Gambling, а найбільші негативні з Income, InflowOutflowsDifference, TotalInflow. Найбільші позитивні значення кореляції отримано між параметрами TotalInflow-Income (0,93), InflowOutflowsDifference- TotalInflow (0,91) та InflowOutflowsDifference-Income (0,84). Найбільші негативні значення кореляції отримано між параметрами InvestmentIncome - CreditMoneyInflows (-0,42), BouncedDirectDebits – InflowOutflowsDifference (-0,25), BouncedDirectDebits – Income (-0,21). Загальний аналіз параметрів кореляції показує відсутність сильної кореляції з Fail, відсутність сильної негативної кореляції між параметрами, а також дуже слабку кореляцію більшості параметрів. Аналіз статистичного розподілу параметрів колонок показав, що для низьких коефіцієнтів кореляції з Fail розподіли величин відрізняються мало. На рис.2 за допомогою box-plot показано відмінності у розподілі величин Gambling та Income для різних Fail у випадку значної позитивної (зліва) та значної негативної (справа) кореляції з Fail. На прикладі параметрів Gambling та Income показано, що значні ігрові ставки та значні доходи діаметрально протилежно впливають на ймовірність отримання кредиту клієнтом.

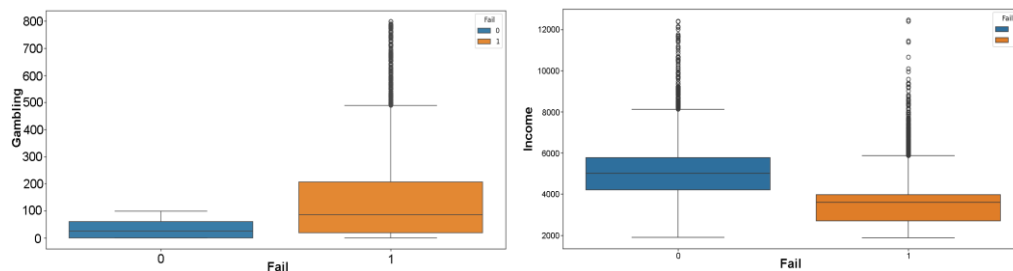


Рис.2. Box-plot для величин Gambling та Income при різних Fail у випадку значної позитивної (зліва) та значної негативної (справа) кореляції з Fail.

Fig. 2. Box-plot for Gambling and Income values at different Fails in case of significant positive (left) and significant negative (right) correlation with Fail.

Побудова ML моделей кредитного рейтингу клієнтів

Класифікаційний аналіз фінансових даних дозволяє значно покращити процеси прийняття рішень у online банківському секторі. Відповідно до [6] високу точність у прогнозуванні кредитного ризику показують класифікатори XGBoost та RandomForest. На цих класифікаторах побудуємо модель, яка здатна класифікувати клієнтів за рівнем кредитного ризику з максимальною точністю. Для навчання та тестування моделей дані розбиті на тренувальний та тестувальний набори у співвідношенні 70% та 30%. Для оцінки моделей використовуємо метрики: точність (accuracy), F1-score, ROC AUC та інші. В таблиці 2 показано метрики класифікації отримані для 3 різних даних. Найкращі метрики класифікації отримано для збалансованих даних 3. Обидва класифікатори XGB та RF показують однакову точність та інші метрики класифікації.

Таблиця 2. Результати класифікації
Table 2. Classification results

Дані	Класифікатори	ROC AUC	Accuracy	Precision		Recall		F1-score	
				0	1	0	1	0	1
Дані 1	XGB	0.75	0.966	0.97	1	1	0.5	0.98	0.67
	RF	0.56	0.943	0.94	1	1	0.12	0.97	0.22
Дані 2	XGB	0.96	0.99	0.99	0.98	1	0.91	1	0.94
	RF	0.95	0.99	0.99	1	1	0.91	1	0.95
Дані 3	XGB	0.98	0.98	0.96	1	1	0.96	0.98	0.98
	RF	0.98	0.98	0.96	1	1	0.96	0.98	0.98

У випадках, коли класи не збалансовані (дані 1 та 2), ROC AUC надає більш надійну оцінку ефективності моделі, ніж проста точність (accuracy). Висока точність accuracy для незбалансованих Даних 1 та 2 пов'язана з тим, що модель з високою точністю передбачає переважаючий клас Fail=0. Однак, модель має низькі значення метрик, таких як чутливість, прецизійність та F1-міра. Це свідчить про те, що модель не може належним чином виявляти менш поширений клас, який відповідає значенню Fail=1. Для даних 3 отримано такі значення метрики ROC AUC від класифікаторів: DecisionTreeClassifier() – 0,97, AdaBoostClassifier() – 0,97, SVC() – 0,86, GaussianNB() – 0,91. Проведено дослідження залежності точності класифікації від кількості компонент за допомогою методу головних компонент PCA(). Значення вектора explained_variance_ratio [0,53 0,38 0,28 0,2 ...] вказують на наявність перших чотирьох важливих компонент, а інші значення вектора на порядок менші. Однак точність класифікації швидко спадає при зменшенні кількості компонент n. Для n=18 і RF класифікатора метрика ROC AUC дорівнює 0,94, для n=15 дорівнює 0,90, для n=12 відповідно 0,79. Проведені дослідження показують, що на основі різноманітної інформації про транзакції клієнтів та їх фінансовий стан, які отримано з OpenBank API можна сформулювати дані для розв'язання проблеми кредитного рейтингу на основі ML класифікаторів. Таким чином сформований датапродукт домену Б забезпечує високу точність класифікатора Домену ML.

Клієнтів розподілено на групи за ймовірністю отримання кредиту. В таблиці 3 показано кількість клієнтів кожної групи в кількісному та відсотковому представленні та запропоновано рекомендації про надання кредиту.

Таблиця 3. Розподіл клієнтів за ймовірностями отримання кредитів
Table 3. Distribution of clients by probabilities of receiving loans

	Fail=0	% користувачів Fail=0	Fail=1	% Fail=1
0.9-1	1212	40,4% (Надати кредит)	17	0,6%
0.7-0.9	221	7,4% (Збільшення відсоткової ставки за кредитом)	32	1,1%
0.5-0.7	56	1,8% (Ручна перевірка даних працівником банку)	10	0,3%
0-0.5	7	0,2% (Відмова. Надання пояснення про причини відмови.)	1445	48,2%

Кількість клієнтів, які неправильно класифіковані моделлю, добре узгоджується з результатами отриманими з confusion matrix. Визначення статистичних параметрів цих клієнтів дозволить підвищити точність класифікації на 1-2%. Інформація про 7 клієнтів (0,2% клієнтів), які не отримали кредит, вказує на можливість перенести їх у групу для ручної перевірки працівником банку. Для всіх клієнтів цієї групи ймовірності отримання кредиту перевищують значення 0,4, а статистичні параметри попадають в діапазон значень клієнтів, які отримали кредит. Для 59 клієнтів в сумі 1,9% спостерігається значний розкид ймовірностей не отримання кредиту від 0,53% до 99%. Клієнтів групи зі ймовірністю отримати кредит більше 0,9 пропонується перевести в групу для ручної перевірки працівником банку.

Висновок

Запропоновано модель data mesh сітки даних з трьома вхідними доменами А, Б, В. Всі домени мають своє операційне джерело даних. Кожна доменна команда будує дата-продукт свого домену в який включає лише очищені, опрацьовані та вибрані дані. Доменні дата-продукти, потім об'єднують в комплексний агрегат, який містить всеохопні дані про сутності в системі. Далі створюють спрямовані на користувача дата-продукти: маркетинговий дата-продукт, дата-продукт ефективності компанії, ML дата продукт. Таким чином data mesh забезпечує децентралізовану та розподілену архітектуру даних для проекту з надання фінансових послуг клієнтам.

В роботі детально проаналізовано створення дата-продуктів доменами: Домен Б та ML домен, а також їх взаємодії. Дата-продукт домену Б будується на основі даних Open Bank API. Open Bank API отримує реальні дані про щоденні детальні транзакції клієнтів, які відповідають інформації, яку клієнт отримує на своїй банківській картці. Дані були категоризовані, агреговані та анонімізовані. Отримано 15 колонок даних за трьома розділами: категоризовані витрати, ризиковані витрати та категоризовані надходження та побудовано дві нові колонки даних, які характеризують різницю доходів та витратків.

Рівень аналізу даних включає домен ML моделей. У домені реалізовано класифікацію даних за допомогою різних класифікаторів. Встановлено, що найвища точність класифікації 0,98 та найвища метрика класифікації ROC AUC 0,98 досягається при використанні класифікаторів XGB та RF на даних отриманих з Домену Б після балансування та аугментації за допомогою Generative Adversarial Network. Результати класифікації та метод головних компонент PCA підтверджують, що побудований дата продукт Домену Б забезпечує високу точності класифікації. Проведено аналіз результатів класифікації. Клієнтів розбито на групи відповідно до ймовірності отримання кредиту. Запропоновано враховувати результати ML аналізу даних для підвищення точності класифікації клієнтів, аналізу фінансових кредитних ризиків та визначення оптимальної відсоткової ставки.

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

- [1] *Dhaouadi A., Bousselmi K., Gammoudi M., Monnet S., Hammoudi S.* Data Warehousing Process Modeling from Classical Approaches to New Trends: Main Features and Comparisons // *Data* **2022**, 7(8), 113. <https://doi.org/10.3390/data7080113>
- [2] *Goedegebuure A., Kumara I., Driessen S., Monsieur G., Tamburri D., Nucci D.* Data Mesh: a Systematic Gray Literature Review // *arXiv:2304.01062v1 [cs.SE]* 3 Apr 2023. <https://doi.org/10.48550/arXiv.2304.01062>

- [3] Hjelkrem L.O., Lange P., Nettet E. The Value of Open Banking Data for Application Credit Scoring: Case Study of a Norwegian Bank // Journal of Risk and Financial Management 15(12):597, 2022. [DOI:10.3390/jrfm15120597](https://doi.org/10.3390/jrfm15120597)
- [4] Shi S., Tse R., Luo W., D'Addona S., Pau G. Machine learning-driven credit risk: a systemic review // Neural Computing and Applications 34(2), 2022. [DOI:10.1007/s00521-022-07472-2](https://doi.org/10.1007/s00521-022-07472-2)
- [5] Strelcenia E., Prakoonwit S. A Survey on GAN Techniques for Data Augmentation to Address the Imbalanced Data Issues in Credit Card Fraud Detection // Machine Learning and Knowledge Extraction 5(1):304-329, 2023. [DOI:10.3390/make5010019](https://doi.org/10.3390/make5010019)
- [6] Khalid A.R., Owoh N., Uthmani O., Ashava M., Osamor J., Adejoh J. Enhancing Credit Card Fraud Detection: An Ensemble Machine Learning Approach// Big Data Cogn. Comput. 2024, 8(1), 6; <https://doi.org/10.3390/bdcc8010006>

DEVELOPMENT OF DATA MESH DATA PLATFORM WITH ML DOMAIN OF DATA ANALYSIS.

M. Fostyak, L. Demkiv

*Ivan Franko National University of Lviv,
50 Drahomanova St., 79005 Lviv, Ukraine
lidia.demkiv@gmail.com*

A data mesh model with three input domains A, B, and C has been proposed. All domains have their own operational data source. Each domain team builds a data product for their domain, which includes only cleaned, processed, and selected data. The domain data products are then combined into a comprehensive aggregate containing all-encompassing data about the entities in the system. Next, consumer-aligned data products are created: a marketing data product, a company performance data product, and an ML data product. Thus, data mesh provides a decentralized and distributed data architecture for a project to deliver financial services to clients.

This study provides a detailed analysis of the creation of data products within Domain B and the ML domain, as well as their interactions. The data product for Domain B is constructed using data from the Open Banking API, which provides real-time data on clients' daily transactions, consistent with the information displayed on their bank statements. The data were categorized, aggregated, and anonymized, resulting in fifteen data columns across three sections: categorized expenditures, risky expenditures, and categorized revenues. Additionally, two new columns were derived to represent the net difference between income and expenditures.

The layer of data analysis includes the ML model domain. In this domain, data classification is implemented using various classifiers. It has been established that the highest classification accuracy of 0.98 and the highest classification metric ROC AUC of 0.98 are achieved when using XGBoost (XGB) and Random Forest (RF) classifiers on data obtained from Domain B after balancing and augmentation with a Generative Adversarial Network. The classification results and the Principal Component Analysis (PCA) method confirm that the data product constructed from Domain B ensures high classification accuracy. A thorough analysis of the classification results was conducted. Clients were segmented into groups based on their probability of obtaining a loan. It is proposed to incorporate the results of ML data analysis to enhance client classification accuracy, analyze financial credit risks, and determine the optimal interest rate.

Keywords: data storage models, data mesh, data domains, classification, ML data analysis.

Стаття надійшла до редакції 22.07.2024.

Прийнята до друку 05.08.2024.