

Вінницький національний технічний університет  
(повне найменування вищого навчального закладу)

Факультет інтелектуальних інформаційних технологій та автоматизації  
(повне найменування інституту, назва факультету (відділення))

Кафедра комп'ютерних наук  
(повна назва кафедри (предметної, циклової комісії))

## МАГІСТЕРСЬКА КВАЛІФІКАЦІЙНА РОБОТА

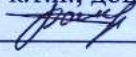
на тему:

«Інформаційна технологія прогнозування успішності  
кінофільму»

Виконав: студент 2-го курсу, групи 2КН-21м  
спеціальності 122 «Комп'ютерні науки»  
(шифр і назва напрямку підготовки, спеціальності)


  
Борисюк В. М.  
(прізвище та ініціали)

Керівник: к.т.н., доцент каф. КН

  
Барабан С. В.  
(прізвище та ініціали)

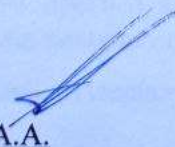
« 15 » 12 2022 р.

Опонент: д.т.н. каф. КСУ

  
Юхимчук М. С.  
(прізвище та ініціали)

« 15 » 12 2022 р.

Допущено до захисту  
Завідувач кафедри КН

  
д.т.н., проф. Яровий А. А.  
(прізвище та ініціали)

« 16 » 12 2022 р.

Вінниця ВНТУ - 2022 рік

Вінницький національний технічний університет  
Факультет інтелектуальних інформаційних технологій та  
автоматизації

Кафедра комп'ютерних наук

Рівень вищої освіти II-й (магістерський)

Галузь знань – 12 “Інформаційні технології”

Спеціальність – 122 “Комп'ютерні науки”

Освітньо-професійна програма – “Системи штучного інтелекту”

**ЗАТВЕРДЖУЮ**

Завідувач кафедри КН

Д.т.н., проф. Яровий А.А.

14. 09

2022 року

## ЗАВДАННЯ

### НА МАГІСТЕРСЬКУ КВАЛІФІКАЦІЙНУ РОБОТУ СТУДЕНТУ

Борисюку Володимиру Миколайовичу

(прізвище, ім'я, по батькові)

1. Тема роботи: Інформаційна технологія прогнозування успішності кінофільму  
керівник роботи к.т.н., доцент кафедри КН Барабан С. В.

затверджені наказом вищого навчального закладу від “14” 09 2022 року № 203

2. Строк подання студентом роботи 18 листопада 2022 року

3. Вихідні дані до роботи:

навчальна вибірка – 3000 кінофільмів, тестова вибірка - 600 кінофільмів; мова програмування – об'єктно-орієнтована.

4. Зміст текстової частини:

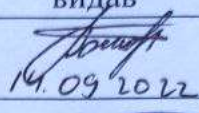

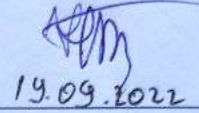
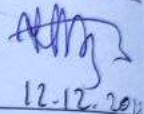
Вступ, аналіз предметної області прогнозування успішності кінофільму, розробка моделі прогнозування успішності кінофільму, розробка інформаційної технології прогнозування успішності кінофільму, програмна реалізація інформаційної технології прогнозування успішності кінофільму, економічна частина, висновки, перелік використаних джерел, додатки.

5. Перелік ілюстративного матеріалу (з точним зазначенням обов'язкових креслень)

Графік кореляції ознак кінофільму до його доходу, графік розподілу доходу кінофільмів з домашньою і без неї, графіки розподілу доходів для фільмів на різних мовах, загальна структура програмного забезпечення прогнозування успішності фільмів, загальна схема бази-даних кінофільмів, схема алгоритму модуля взаємодії з базою даних, схема алгоритму модуля імпортування фільмів з

tmdb, схема алгоритму модуля прогнозування успішності кінофільму, вікно програми з результатом прогнозування успішності кінофільму.

6. Консультанти розділів роботи

Розділ	Прізвище, ініціали та посада консультанта	Підпис, дата	
		завдання видав	виконання прийняв
1-4	Барабан С. В., к.т.н., доц. каф. КН	 14.09.2022	 14.09.2022
5	Буреннікова Н. В., д. е. н., проф. каф. ЕПВМ	 19.09.2022	 12.12.2022

7. Дата видачі завдання 14.09 2022 року

КАЛЕНДАРНИЙ ПЛАН


№ з/п	Назва етапів магістерської кваліфікаційної роботи	Строк виконання етапів роботи
1	Аналіз сучасних підходів прогнозування успішності кінофільмів.	12.08.2022 - 01.09.2022
2	Побудова моделей прогнозування успішності кінофільму.	01.09.2022 - 15.09.2022
3	Практичне застосування та оцінка ефективності розроблених моделей	15.09.2022 - 20.11.2022
4	Підготовка економічної частини	20.11.2022 - 29.11.2022
5	Апробація та/або впровадження результатів дослідження	29.11.2022 - 10.12.2022
6	Оформлення пояснювальної записки, графічного матеріалу та презентації	10.12.2022 - 16.12.2022

Студент

Керівник роботи

  
(підпис)

Борисюк В. М.

  
(підпис)

Барабан С. В.

## АНОТАЦІЯ

УДК 004.8

Борисюк В. М. Інформаційна технологія прогнозування успішності кінофільму. Магістерська кваліфікаційна робота зі спеціальності 122 – комп'ютерні науки, освітня програма - комп'ютерні науки. Вінниця: ВНТУ, 2022. 103 с.

На укр. мові. Бібліогр.: 28 назв; рис. 40; табл. 26.

У даній магістерській кваліфікаційній роботі на основі проведеного аналізу розроблено програмного забезпечення для прогнозування успішності кінофільму, використовуючи нейромережевий підхід. Запропоновано власну структуру нейронної мережі. Розроблено інформаційну технологія прогнозування успішності кінофільму. Реалізоване програмне забезпечення для прогнозування успішності кінофільму. Серверна частина реалізована на мові програмування Python у програмному середовищі PyCharm, клієнтська частина реалізована у вигляді додатку на мові програмування Dart, з використанням фреймворку Flutter. Розроблений продукт характеризується відсутністю прямих аналогів, наявністю графічного інтерфейсу користувача та підвищеною точністю прогнозування успішності кінофільму.

Графічна частина складається з 7 плакатів.

У економічному розділі при оцінюванні за технічними параметрами, згідно узагальненого коефіцієнту якості розробки, визначено, що науково-технічна розробка переважає існуючі аналоги приблизно в 1,56 рази. Також термін окупності становить 0,96 р., що менше 3-х років, що свідчить про комерційну привабливість науково-технічної розробки

Ключові слова: прогнозування, штучний інтелект, нейронні мережі, регресія, *python*.

## **ABSTRACT**

Borysiuk V. M. Information technology for the signature dynamic invariant features formation. Master's thesis in the speciality 122 - Computer Sciences, educational program - Computer Sciences. Vinnytsia: VNTU, 2022. 103 p.

In Ukrainian language. Bibliogr.: 28 titles; fig. 40; table 26.

In this master's thesis, based on the analysis, software was developed for predicting the success of a movie using a neural network approach. A proprietary structure of a neural network is proposed. An information technology for predicting the success of a motion picture has been developed. Implemented software for predicting the success of a motion picture. The server part is implemented in the Python programming language in the PyCharm programming environment, the client part is implemented as an application in the Dart programming language, using the Flutter framework. The developed product has the following advantages: the absence of direct analogues, the presence of a graphical user interface and increased accuracy in predicting the success of a motion picture.

The graphic part consists of 7 posters.

In the economic section, when evaluating according to technical parameters, according to the generalized coefficient of development quality, it was determined that scientific and technical development prevails over existing analogues by approximately 1.86 times. Also, the payback period is 0.96 years, which is less than 3 years, which indicates the commercial attractiveness of scientific and technical development.

**Keywords:** prediction, artificial intelligence, neural networks, regression, python.

## ЗМІСТ

ВСТУП .....	7
1 АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ .....	10
1.1 Постановка задачі прогнозування успішності кінофільму .....	10
1.2 Огляд відомих методів прогнозування успішності кінофільму .....	11
1.3 Аналіз об'єкту проектування .....	18
1.4 Висновок до розділу 1 .....	18
2 РОЗРОБКА МОДЕЛІ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ .....	19
2.1 Виділення ознак кінофільмів для прогнозування успішності фільму ...	19
2.1.1 Відношення бюджету кінофільму до доходу .....	23
2.1.2 Зв'язок між домашньою сторінкою та доходом. ....	25
2.1.3 Зв'язок між оригінальною мовою фільму (original_language) і середнім доходом. ....	26
2.1.4 Найбільш вживані слова в фільмах. ....	27
2.1.5 Вплив опису фільму на прибуток .....	29
2.1.6 Вплив дати виходу на прибуток кінофільму .....	31
2.1.7 Співвідношення між хронометражем кінофільму та доходом. ....	34
2.2 Розробка моделі прогнозування успішності кінофільму .....	34
2.2.1 Метод лінійної регресії .....	35
2.2.2 Метод Random Forest. ....	36
2.2.3 Метод Gradient Boosting. ....	37
2.2.4 Метод регресії на основі нейронної мережі. ....	37
2.3 Висновок до розділу 2 .....	39
3 РОЗРОБКА ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ .....	40
3.1 Розробка структури програмного забезпечення прогнозування успішності кінофільму .....	40

3.2	Проектування бази-даних інформаційної технології прогнозування успішності фільмів.....	42
3.3	Розробка алгоритму модуля взаємодії з базою даних кінофільмів.....	46
3.4	Розробка алгоритму модуля імпортування кінофільмів з TMDb.....	48
3.5	Розробка алгоритму модуля редагування записів бази даних фільмів..	49
3.6	Розробка алгоритму модуля прогнозування успішності кінофільму ....	50
3.7	Висновок до розділу 3.....	51
4	ПРОГРАМНА РЕАЛІЗАЦІЯ ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ .....	52
4.1	Обґрунтування вибору мови програмування .....	52
4.2	Обґрунтування вибору середовища програмування .....	53
4.3	Програмна реалізація модулів інформаційної технології прогнозування успішності кінофільмів.....	55
4.3.1	Програмна реалізація модуля взаємодію з базою даних.....	55
4.3.2	Програмна реалізація модуля редагування записів в базі даних .....	57
4.3.3	Програмна реалізація модуля імпортування кінофільмів з TMDb ..	58
4.3.4	Програмна реалізація модуля прогнозування успішності кінофільму.....	59
4.4	Тестування та аналіз результатів роботи програми прогнозування успішності кінофільму.....	61
4.5	Висновок до розділу 4.....	64
5	ЕКОНОМІЧНА ЧАСТИНА .....	65
5.1	Проведення комерційного та технологічного аудиту науково-технічної розробки .....	65
5.2	Розрахунок узагальненого коефіцієнта якості розробки .....	70
5.3	Розрахунок витрат на проведення науково-дослідної роботи.....	72

5.3.1	Витрати на оплату праці.....	72
5.3.2	Відрахування на соціальні заходи .....	75
5.3.3	Сировина та матеріали.....	75
5.3.4	Розрахунок витрат на комплектуючі.....	76
5.3.5	Спецустаткування для наукових (експериментальних) робіт .....	76
5.3.6	Програмне забезпечення для наукових (експериментальних) робіт	77
5.3.7	Амортизація обладнання, програмних засобів та приміщень .....	78
5.3.8	Паливо та енергія для науково-виробничих цілей.....	79
5.3.9	Службові відрядження.....	80
5.3.10	Витрати на роботи, які виконують сторонні підприємства, установи і організації.....	81
5.3.11	Інші витрати.....	81
5.3.12	Накладні (загальновиробничі) витрати.....	81
5.4	Розрахунок економічної ефективності науково-технічної розробки при її можливій комерціалізації потенційним інвестором .....	83
5.4	Висновок до розділу 5.....	87
	ВИСНОВКИ.....	88
	ПЕРЕЛІК ВИКОРИСТАНИХ ДЖЕРЕЛ.....	90
	Додаток А (обов'язковий) Результат перевірки на плагіат в онлайн-системі UNICHECK .....	94
	Додаток Б (обов'язковий) Лістинг програми .....	95
	Додаток В (обов'язковий) ІЛЮСТРАТИВНА ЧАСТИНА.....	<b>Ошибка!</b>
	<b>Закладка не определена.</b>	
	Додаток Г (довідниковий) Інструкція користувача .....	108



## ВСТУП

**Актуальність.** Кіноіндустрія відіграє дуже важливу роль у нашому житті, вона має вирішальний вплив на суспільство, розширює наші знання про історію та культуру, надихає на розважає. Фільми є високопрофесійною багатомільярдною індустрією з сотнями нових фільми, створені щороку.

Основні складності процесу передбачення успішності кінофільму:

- визначити методи попередньої обробки наборів даних;
- які ознаки є найбільше корисними;
- оцінка показників успішності моделі передбачення.

Для того щоб вирішувати складні і погано формалізовані завдання і виник напрямок, який називається штучні нейронні мережі. Штучні нейронні мережі складаються з нейроноподібних елементів, з'єднаних між собою в мережу. Нейронні мережі знайшли застосування практично.

Якщо можна за допомогою комп'ютера передбачити, наскільки вдалим буде фільм, навіть до його релізу, це був би потужний інструмент для використання. Починаючий голлівудський режисер або кіностудія з деякими технічними навичками могли б передбачити, чи їхня ідея стане успішною і чи це є безпечною інвестицією.

**Зв'язок роботи з науковими програмами, планами, темами.** Магістерська робота виконана відповідно до напрямку наукових досліджень кафедри комп'ютерних наук Вінницького національного технічного університету 22 К1 «Моделі, методи, технології та пристрої інтелектуальних інформаційних систем управління, економіки, навчання та комунікацій» та плану наукової та навчально-методичної роботи кафедри.

**Мета і завдання досліджень.** Метою магістерської кваліфікаційної роботи є підвищення точності прогнозування успішності кінофільму програмними засобами за рахунок використання штучних нейронних мереж.

Для досягнення мети розробки необхідно виконати такі задачі:

- провести аналіз предметної області передбачення успішності кінофільмів;
- розглянути існуючі методи передбачення успішності кінофільмів та обрати й обґрунтувати вибір методу, який задовольняє мету даної магістерської кваліфікаційної роботи;
- розробити математичну модель передбачення успішності кінофільмів;
- сформулювати стадії інформаційної технології, розробити структуру та алгоритм роботи програмного засобу;
- виконати програмну реалізацію запропонованої інформаційної технології;
- провести тестування програмного продукту та виконати аналіз отриманих результатів.

**Об’єкт дослідження** – процес прогнозування успішності кінофільмів комп’ютерними засобами з використанням нейронних мереж.

**Предмет дослідження** – інформаційна технологія та програмні засоби прогнозування успішності кінофільмів з використанням нейронних мереж та достовірність їх роботи.

**Методи дослідження.** У роботі використані наступні методи наукових досліджень: системного аналізу, теорії штучних нейронних мереж для реалізації інформаційної технології, методи математичної статистики для розробки процесу розв’язання задачі нейромережевого передбачення успішності кінофільмів та обрахунків результатів експериментів із програмним засобом, об’єктно-орієнтованого програмування.

**Наукова новизна одержаних результатів.**

Набула подальшого розвитку інформаційна технологія передбачення успішності кінофільмів, в якій на відміну від існуючих використовується нейронна мережа, що дозволяє підвищити точність прогнозування.

**Практичне значення** одержаних результатів полягає в тому, що на основі проведених досліджень розроблено програмне забезпечення для передбачення успішності кінофільмів.

Запропонована інформаційна технологія сприяє підвищенню точності програмних засобів передбачення успішності кінофільмів, зокрема:

- розроблено алгоритм роботи програмного передбачення успішності кінофільмів;
- розроблено програмні засоби для передбачення успішності кінофільмів.

**Достовірність теоретичних положень** магістерської кваліфікаційної роботи підтверджується коректністю постановки завдання, коректністю використання математичного апарату методів дослідження, експериментальними дослідженнями, тестування програмної реалізації інформаційної технології розпізнавання музичних патерів. Адекватність розроблених математичних моделей підтверджується результатами експериментальних досліджень.

**Особистий внесок здобувача.** Усі результати, наведені у магістерській кваліфікаційній роботі, отримані самостійно. У працях, написаних у співавторстві, здобувачу належать: виділення ознак кінофільму, що впливають на його успішність [1].

**Апробація результатів роботи.** Результати досліджень апробовані на конференції «Молодь в науці: дослідження, проблеми, перспективи-2023», Вінниця, 15 листопада 2022 – 12 травня 2023 року.

**Публікації.** За результатами досліджень опубліковано одні тези доповіді на науково-технічній конференції [1] та подано заяву на авторське право на програмне забезпечення.

# 1 АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ

## 1.1 Постановка задачі прогнозування успішності кінофільму

Прогнозування рейтингу фільму може допомогти людям визначити, який фільм потрібно дивитися. Експерти кіноіндустрії сходяться на думці, що прибуток є ключовим фактором успіху кіно і допомагає кіно виробникам та інвесторам досягти фінансового успіху [2]. Компанії можуть побачити, які фільми, ймовірно, мають хороші рейтинги, і розробити стратегії розробки фільму для збільшення прибутку, наприклад, випускаючи товари для фільмів або створюючи події та рекламні акції, пов'язані з фільмом.

Крім того, використовуючи історичні значення, отримані з раніше випущених фільмів, можна прогнозувати прибуток до того, як фільм буде знятий. Щоб уникнути втрат, компанії-виробники фільмів можуть складати стратегічні плани та приймати рішення щодо виходу фільму в прокат.

Прогнозування касових зборів конкретного фільму привернуло багатьох вчених, тому що це передбачення є важкою та складною проблемою. Часто кіноіндустрія залишає у людей враження прибуткової сфери. Цьому враженню сприяють образи знаменитостей та валовий прибуток, який вимірюється сотнями мільйонів доларів. Однак більшість людей лише звертають увагу до найуспішніших фільмів, які, як правило, приносять певний прибуток, але в цілому це враження невірне. За статистикою, з будь-яких десяти великих кінофільмів, створених у середньому, може бути шість-сім кінофільмів характеризується як збитковий [3]. Ці цифри свідчать що кіноіндустрія є одним із найризикованіших ринків індустрії розваг, що виправдовує високу віддачу від успішних фільмів. Це через ці високі ризики у створенні фільмів, які складають адекватний і точний бюджетний план прогнозування доходів стає дуже важливим. Хоча є кілька досліджень щодо прогнозування рейтингу фільмів. Ефективність прогнозування рейтингу перед виходом фільму в

прокат все ще потрібно підвищити. Це дослідження зосереджено на прогнозуванні рейтингу на основі даних, які можна отримати до виходу фільму в прокат. Атрибути метаданих, отримані з IMDb і TMDb, були використані як предикатори. Для створення функцій використовуються такі атрибути, як сюжет, артисти, персонал, режисер, жанр, рейтинг вмісту та рейтинг.

Однією з унікальних рис, які використовуються в цьому дослідженні, є історичні цінності. Ці функції були створені на основі зв'язку між фільмом і раніше випущеними фільмами. Передбачається, що прогнозований рейтинг, заснований на цих історичних значеннях, безумовно, є більш об'єктивним, ніж рейтинг від аудиторії, який з'явився, коли фільм щойно вийшов у прокат. Метод когортного оцінювання шукає фільми, які мають схожість на основі наявних історичних атрибутів і значень, і робить прогноз касового прибутку на основі цих подібностей.

## 1.2 Огляд відомих методів прогнозування успішності кінофільму

**Модель дифузії Басса.** Однією з перших моделей поведінки споживачів, які можна застосувати для моделювання касових зборів фільму, є модель Баса [4]. У моделі Басса використовується звичайний диференціал

Рівняння, яке пов'язує, як нові продукти вживаються в популяції:

$$\frac{f(t)}{1 - F(t)} = p + qF(t), \quad (1.1)$$

де  $f(t)$  – швидкість зміни встановленої базової фракції;  $F(t)$  – встановлена базова частка;  $p$  – коефіцієнт інноваційності;  $q$  – коефіцієнт обмеження.

Розв'язком цієї системи дає відому S-криву (Рисунок 1.1), яка

представляє швидкість прийняття продукту з часом. Ця модель була однією з найвпливовіших і поширених цитованих робіт з історії науки управління.

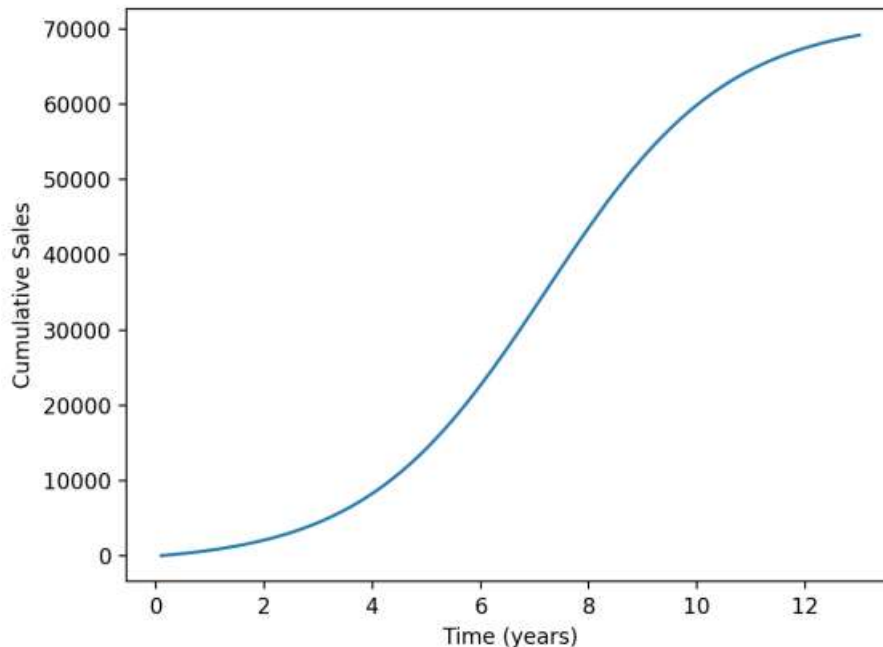


Рисунок 1.1 - S-крива, що показує загальні продажі за період часу

У роботі «*A Parsimonious Model for Forecasting Gross Box-Office Revenues of Motion Pictures*» [5] представлена модель, що використовує експоненційний розподіл з двома параметрами та три параметри гамма-розподілу, яка моделює процес, за допомогою якого людина вирішить подивитися фільм у два етапи. По-перше, який час, необхідний для прийняття рішення про те, чи варто чи ні подивитися фільм. По-друге, змоделювати час, необхідний для виконання цього рішення. У роботі «*The motion picture industry: critical issues in practice, current research, and new research directions*» [6] розробили програму для прогнозування касових зборів для фільму, використовуючи лише дані перед релізом. Зокрема, дана програма здатна передбачити касові збори для фільму на основі моделювання поведінки споживачів за допомогою інтерактивних ланцюгів Маркова [7], де ймовірності переходу залежать від кількості людей, які вже перебувають в інших регіонах.

**Модель Едвардса-Бакмайра.** Едвардс і Бакмайр розробили систему поєднаних звичайних диференціальних рівняння для моделювання касових зборів фільму після випуску [8]. Модель Едвардса-Бакмайра використовує детермінований підхід, використовуючи керуючі рівняння зі швидкістю зміни доходу в момент  $t$ , спочатку заданою:

$$\frac{d\tilde{G}}{dt} = \tilde{S}\tilde{A}, \quad (1.2)$$

$$\frac{d\tilde{A}}{dt} = -a_A\tilde{A}, \quad (1.3)$$

$$\frac{d\tilde{S}}{dt} = -a_S\tilde{S} \quad (1.4)$$

де,  $\tilde{G}$ ,  $\tilde{A}$  та  $\tilde{S}$  - прибуток, кількість зароблених грошей за показ за тиждень і кількість показів на якому представлений фільм, відповідно.  $P$  - це ціна квитка. Умови включають  $\tilde{G}(0) = 0$  та  $\tilde{G}(\infty) = \int_0^\infty \tilde{S}\tilde{A}dt$ .

Модель Едвардса-Бакмайра розвивається далі, намагаючись моделювати негативну реакцію людини ( $\tilde{H}$ , с відсоток людей, які ненавидять фільм, наданий  $H_0$ ) до фільму, яким керує:

$$\frac{d\tilde{H}}{dt} = \frac{H_0\tilde{S}\tilde{A}}{P} \quad (1.5)$$

Це передбачає попередні знання про загальну кількість людей, які будуть дивитися фільм, і знання кількості людей, яким він не подобається, що недоцільно в реальному світі. Модель Едвардса-Бакмайра додатково розвивається, дозволяючи людям дивитися фільм кілька разів, додаючи параметри для врахування жанру, суми, витраченої на рекламу, та ефективності реклами. Більшість з цих параметрів необхідно оцінювати в

режимі реального часу.

Модель Едвардса-Бакмайра дала багатообіцяючі результати для ряду фільмів, що зображено на рисунках 1.2, 1.3 та 1.4.

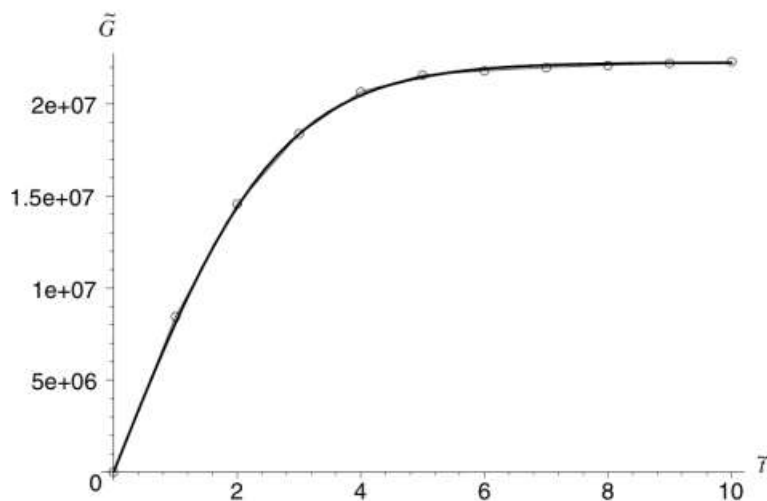


Рисунок 1.2 - Прибуток для фільму *At First Sight*. Крива передбачення

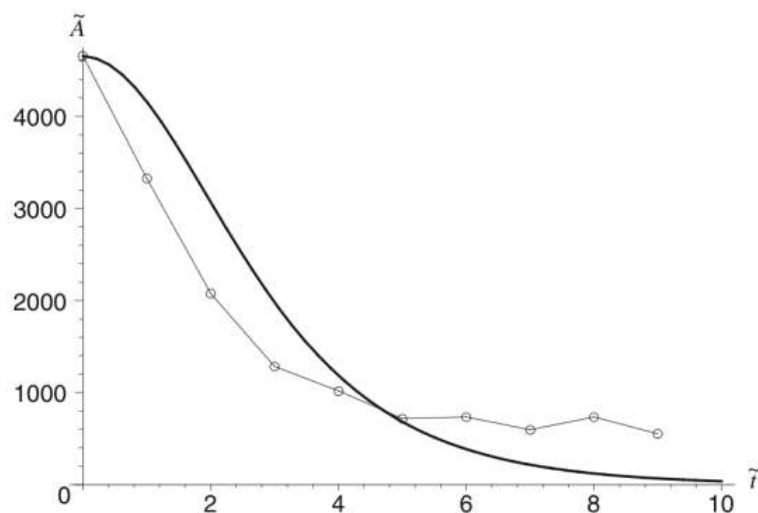


Рисунок 1.3 - Аудиторія для фільму *At First Sight*. Крива передбачення



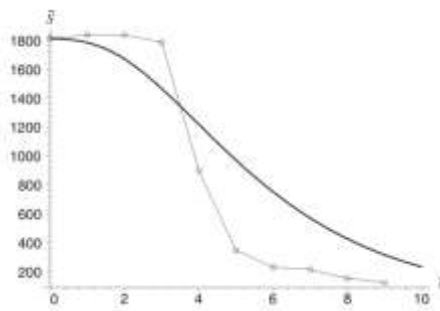


Рисунок 1.4 – Кількісит показів для фільму *At First Sight*. Крива передбачення

**Методи стохастичної апроксимації.** Касові збори фільмів — це випадковий процес. Прибуток фільму визначається кількістю кінотеатрів, які демонструють фільм, і кількість людей, які йдуть дивитися його.

У цьому є випадковість, і стохастичні моделі намагаються змоделювати цю випадковість.

Стохастичні процеси можна моделювати за допомогою ланцюгів Маркова, де наступний крок залежить лише від поточного стану. Марківський процес безпам'ятний – він не знає, як було досягнуто поточного стану або в який час він був досягнутий. Наприклад, за  $n$  станів  $S = \{s_1, s_2, \dots, s_n\}$ , марковський процес може переходити з одного стану в інший з ймовірністю  $P_{ij}$  (у процесі дискретного часу). На рисунку 1.5 зображено ланцюг Маркова.

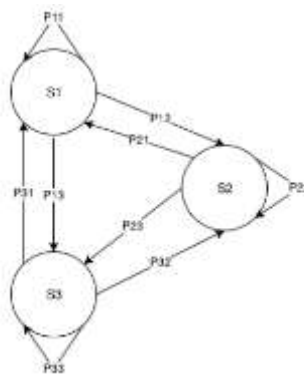


Рисунок 1.5 – Ланцюг Маркова з трьома станами

Ланцюги Маркова може моделювати поведінку доходів окремих

кінотеатрів, але перевести це на макроскопічний погляд на всю систему може бути складно. Моделювання кількох людей за допомогою ланцюгів Маркова призводить до проблеми вибуху простору стану [9], що може зробити цей підхід нежиттєздатним.

Підходячи до цього з детермінованої точки зору, моделюючи стохастичний процес за допомогою звичайних диференціальних рівнянь. Це означає, що ми втрачаємо дискретні стани моделі ланцюга Маркова, але отримуємо модель, яку набагато легше моделювати, ніж кілька процесів Маркова. Однак ігнорування стохастичних ефектів може означати, що в деяких випадках рівняння створюють значну помилку.

**Моделі машинного навчання.** Машинне навчання (ML) — це підмножина галузі штучного інтелекту (AI), яка розробляє алгоритми, які можуть навчатися на даних. Контрольований ML включає модель, яка вивчає функцію з позначених навчальних даних.

У разі прогнозування продуктивності фільму історичні дані, що містять таку інформацію, як бюджет, мова тощо, можуть бути використані для навчання моделі. Функція Інженерія даних є важливою частиною розробки ефективної моделі ML. Конструювання ознак даних а є важливою частиною розробка ефективної моделі ML.

Поставлена тут проблема полягає в регресії, де вихідні дані будуть числом з дійсним значенням, наприклад, загальний касовий збір фільму.

Прогнозування касових зборів з певним успіхом розглядалося як проблема класифікації і є потенційно розумний підхід до проблеми, якщо регресія не може дати хороших результатів.

**Лінійна регресія.** Лінійний регресійний аналіз – це підхід до моделювання зв'язку між незалежними змінними та залежними змінними.

Якщо похідну неможливо обчислити, то ітераційне рішення можна знайти за допомогою градієнтного спуска або при використанні великого набору даних за допомогою стохастичного градієнтного спуску [10].

**Нейронні мережі.** Нейронна мережа (NN) — це система, натхненна структурою біологічних нейронних мереж. NN можуть бути структурованими так, щоб мати кілька шарів, кожен нейрон з'єднаний з кожним нейроном у наступному шару (крім вхідного шару). На рисунку 1.6 зображено схему нейронної мережі.

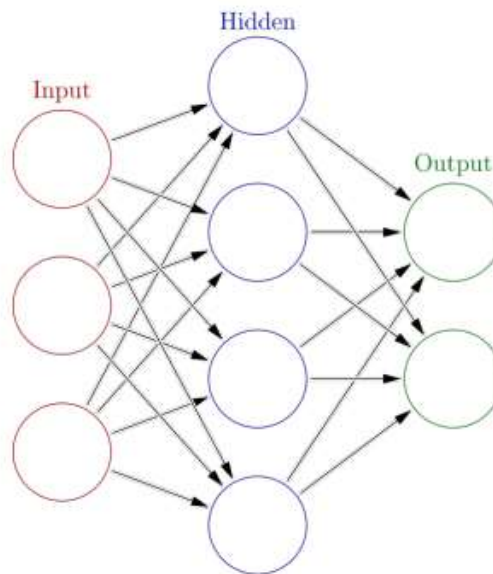


Рисунок 1.6 – Ілюстрація тришарової нейронної мережі

Кожен нейрон пов'язаний із вагою, який нейрон використовує разом зі своїми вхідними параметрами для обчислення вихідного значення. Функцію активації можна застосувати до суми зважених вхідних даних з попереднього шару. Застосована функція активації є гіперпараметром, який необхідно вибрати. Нейронні мережі підлаштовуються до навчального набору даних за допомогою алгоритму зворотного поширення, який може використовувати алгоритми оптимізації, такі як стохастичний градієнтний спуск, які надають додаткові параметри для оптимізації, включаючи швидкість навчання та розмір пакету.

### **1.3 Аналіз об'єкту проектування**

Основною метою розробки інформаційної технології прогнозування успішності кінофільмів є надання кінокомпаніям інформації про успішність кінофільму, що ще знаходиться в розробці. Це дозволить компаніям скорегувати напрямки своїх фільмів, щоб отримати максимальну віддачу від аудиторії.

Вимоги інформаційної технології:

- Зручний інтерфейс користувача.
- Наявність персонального кабінету.
- Можливість імпортувати власні дані про фільми.
- Актуальність – система повинна весь час наповнюватись актуальною інформацією, для передбачень.
- Швидкість надання передбачення.
- Інтерпритованість передбачень.

### **1.4 Висновок до розділу 1**

У даному розділі було проведено аналіз предметної області прогнозування успішності кінофільму, а саме – сформульовано постановку задачі, проведено огляд відомих методів розв'язання задачі прогнозування успішності кінофільму.

Виходячи з аналізів предметної області, дало можливість визначити основні вимоги до системи прогнозування успішності кінофільму.

## 2 РОЗРОБКА МОДЕЛІ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ

Методологія розробки інформаційної технології прогнозування успішності кінофільму:

- збір даних з набору даних *TMDB*;
- дослідницький аналіз даних, конструювання ознак, візуалізація даних, дослідження кореляції ознак;
- моделювання експериментів для оцінки продуктивності та вибору методу машинного навчання;
- оцінка моделі на перевірочному наборі даних.

Усі дослідження проводяться з використанням наступних інструментів:

- Python;
- Jupyter;
- Pandas;
- Numpy;
- Scikit-learn, wordcount, elit, TFID.

### 2.1 Відокремлювання ознак кінофільмів для прогнозування успішності фільму

Для знаходження ознак кінофільм, що впливають на успішність кінофільму використано набір даних *TMDB* [11]. *TMDB* – один із найбільших база даних про фільми, що містить інформацію про 568 729 кінофільмів та серіалів. На рисунку 2.1 зображено прикладу набору даних *TMDB*.

id	belongs_to_collection	budget	genres	homepage	imdb_id	original_language	original_title
0	1 [ {'id': 315578, 'name': 'Hot Tub Time N...'} ]	34000000	[ {'id': 35, 'name': 'Comedy'} ]		tt2057294	en	Hot
1	2 [ {'id': 107674, 'name': 'The Princess D...'} ]	40000000	[ {'id': 35, 'name': 'Comedy'}, {'id': 1...} ]		tt0358931	en	The Princess Diaries
2	1	2300000	[ {'id': 38, 'name': 'Drama'} ]	http://www.fox.com/movies/	tt2982860	en	
3	4	1200000	[ {'id': 33, 'name': 'Thriller'}, {'id': ...} ]	http://www.fox.com/	tt1821680	en	
4	5	0	[ {'id': 28, 'name': 'Action'}, {'id': 8...} ]	http://www.fox.com/	tt1388112	en	

Рисунок 2.1 – приклад набору даних *TMDB*

В таблиці 2.1 перелічено усі стовбці даних та їх тип даних.

Таблиця 2.1 – Стовбці набору даних

Колонка	Тип даних	Опис
id	integer	Індикатор фільму
belongs_to_collection	json	Приналежність до франшизи
budget	integer	Бюджет
genres	json	Список жанрів фільму
homepage	string	Посилання на домашню сторінку
imdb_id	string	Індифікатор фільму на IMDb
original_language	string	Оригінальна мова фільму
original_title	string	Заголовок фільму
overview	string	Короткий опис фільму
popularity	float	Рейтинг популярності
poster_path	string	Посилання на постер
production_companies	json	Список компаній виробників

Продовження таблиці 2.1

Колонка	Тип даних	Опис
production_countries	json	Список країн виробників
release_date	date	Дата релізу
runtime	integer	Хронометраж в хвилинах
status	stiring	Статус
keywords	json	Ключові слова
cast	json	Список акторів
crew	json	Знімальна команда
revenue	integer	Прибуток

**Кореляція ознак.** Діапазон значень коефіцієнта кореляції становить від -1.0 до 1.0. Якщо кореляція дорівнює нулю, це означає, що між двома змінними немає лінійного зв'язку. Якщо кореляція наближається до 1, це означає що дві зміни сильно пов'язані між собою. Наприклад, ціна на нафту безпосередньо пов'язана з цінами на авіаквитки. Якщо значення кореляції наближається до -1, це означає що зміни негативно корельовані, тобто якщо одна зміна збільшується, інша зміна зменшується з таким самим співвідношенням і навпаки. На рисунку 2.2 зображено кореляцію ознак кінофільмів.



Рисунок 2.2 – Кореляція ознак кінофільмів

З рисунка 2.2 видно що кореляція доходу за бюджетом становить 0.75 одиниць, а кореляція між доходом і хронометражем фільму становить 0.22 одиниць. Це пояснює, що прибуток тісно пов'язаний із бюджетом фільму, також з доходом пов'язаний тривалість фільму.

На рисунку 2.3 зображено графік розподіл доходу.



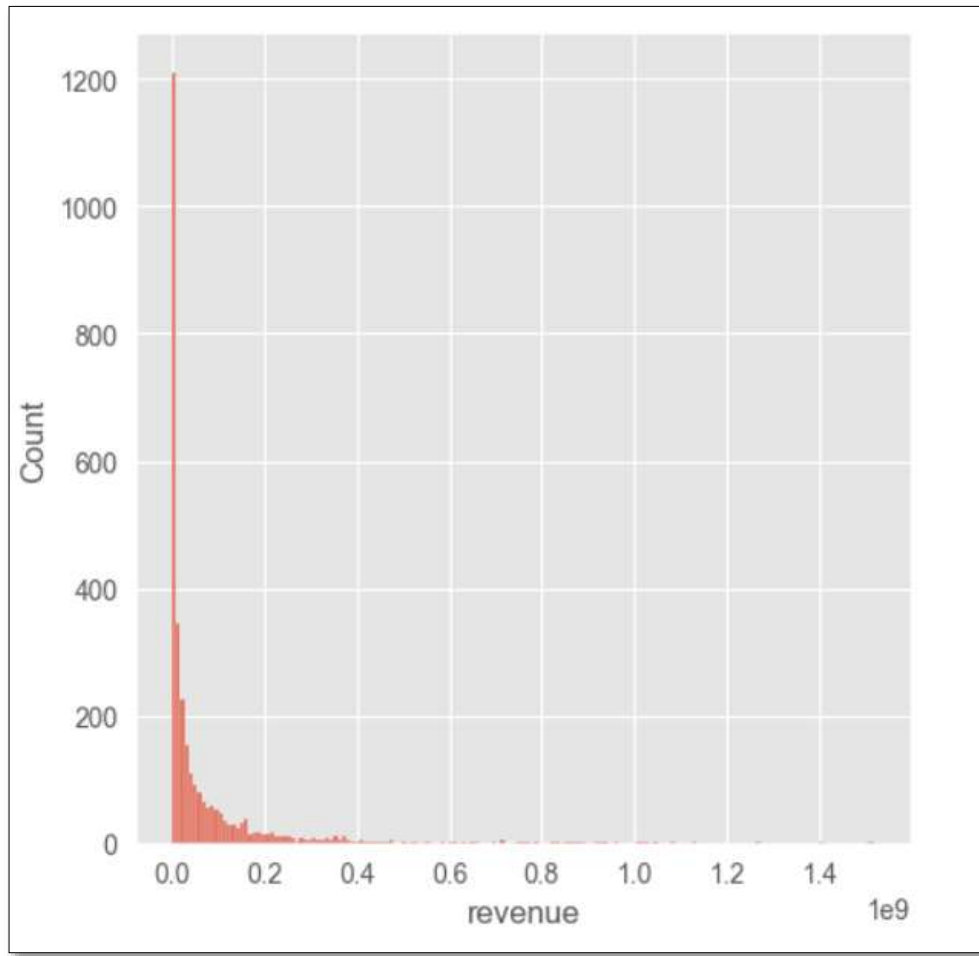


Рисунок 2.3 – Графік розподіл доходу

З рисунку 2.3 видно, що дані дуже спотворені, тому важко зробити висновки з графіку, тому потрібно нормалізувати дані. Для цього використовується логарифмічне трансформування, формула (1).

$$\log_1 p = \log(1 + x) \quad (2.1)$$

#### 2.1.1 Відношення бюджету кінофільму до доходу

На рисунку 2.4 зображено результат логарифмічного трансформування. Після трансформування дані стали нормально розподілені, що має меншу асиметрію та ексцес.

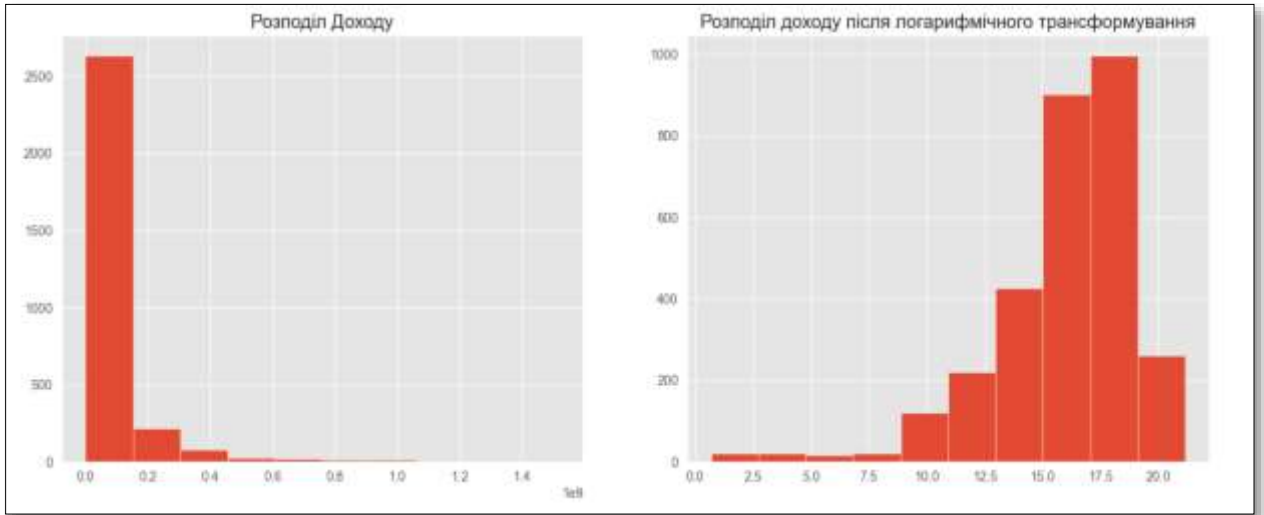


Рисунок 2.4 – Логарифмічне трансформування графіку доходу

На рисунку 2.5 зображено результат логарифмічного трансформування для бюджету.

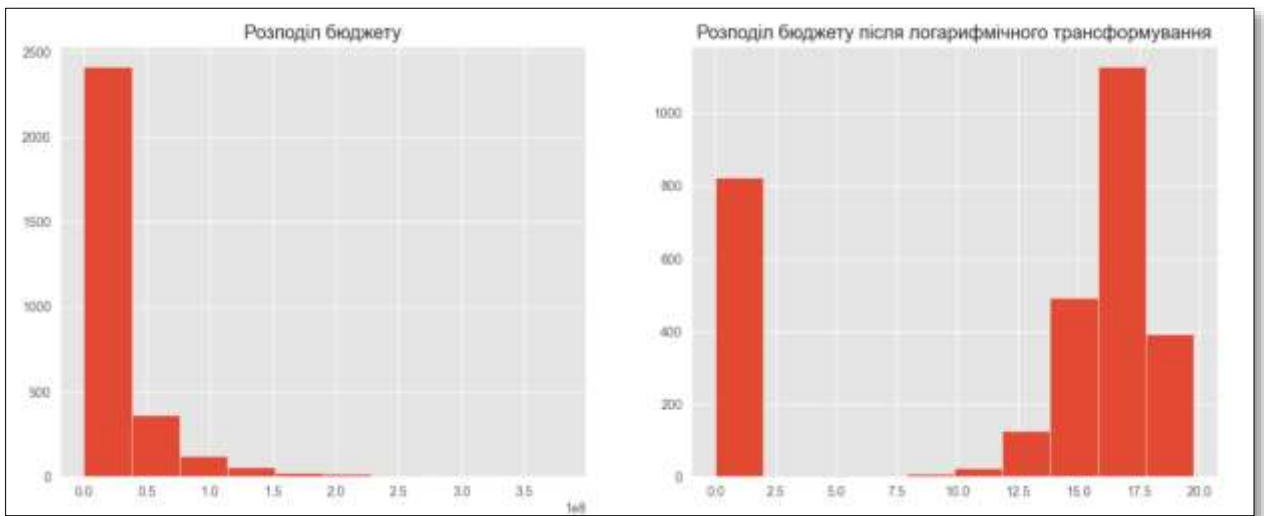


Рисунок 2.5 – розподіл бюджету

Наступний кроком в дослідженні даних буде побудова діаграми розсіювання доходу.

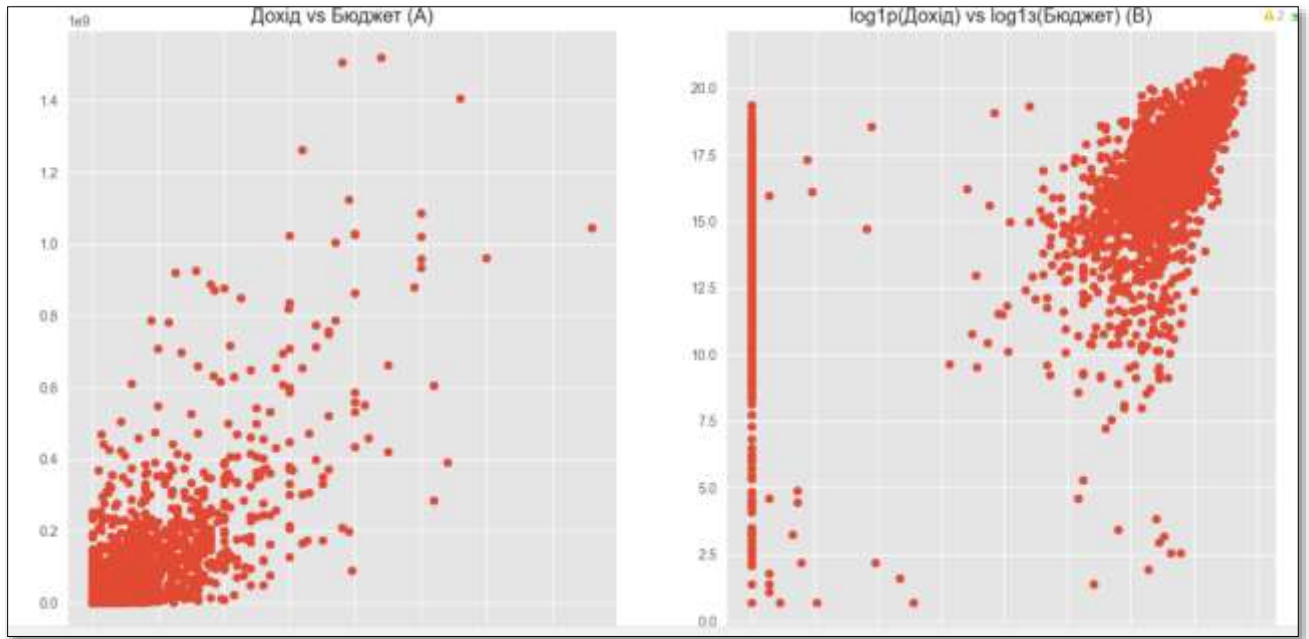


Рисунок 2.6 – Графік діаграми розсіювання бюджету проти доходу

На графіку 2.6 (A) видно, що бюджет та прибуток має певну міру кореляції. Також можна побачити, що в наборі даних є 815 фільмів з нульовим бюджетом.

### 2.1.2 Зв'язок між домашньою сторінкою та доходом.

Також важливою ознакою успішності кінофільму є наявність домашньої сторінки фільму, де користувачі можуть дізнатись усю необхідну інформацію про кінофільм та слідкувати за оновленнями. На рисунку 2.7 зображено статистику доходу фільмів з домашньою сторінкою та без неї.

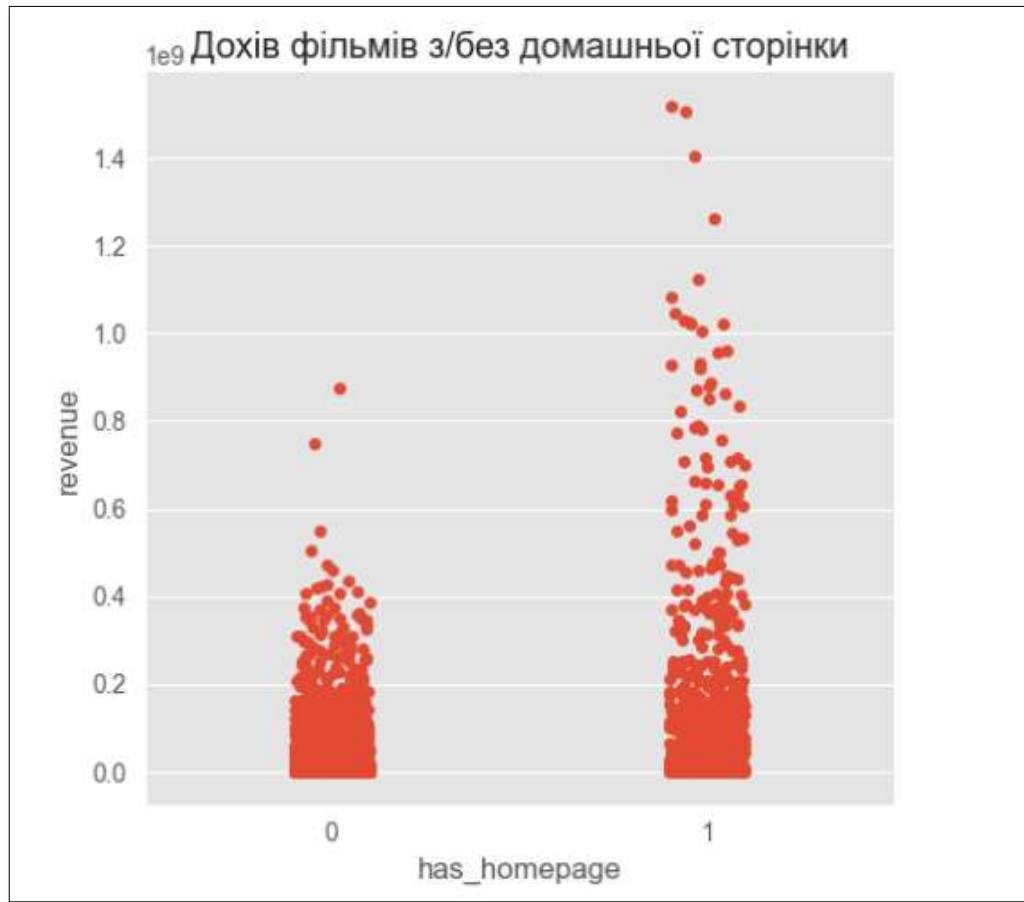


Рисунок 2.7 – Графік доходу фільмів з домашньою сторінкою і без неї

З рисунка 2.7 видно, що фільми, які мають домашню сторінку, мають більший прибуток порівняно з фільмами, що не мають домашньої сторінки. Отже, звідси можна зробити висновок, що наявність домашньої сторінки та прибуток корелюють між собою.

2.1.3 Зв'язок між оригінальною мовою фільму (`original_language`) і середнім доходом.

На рисунку 2.8 зображено графік розподілу доходів між різними мовами фільмів.

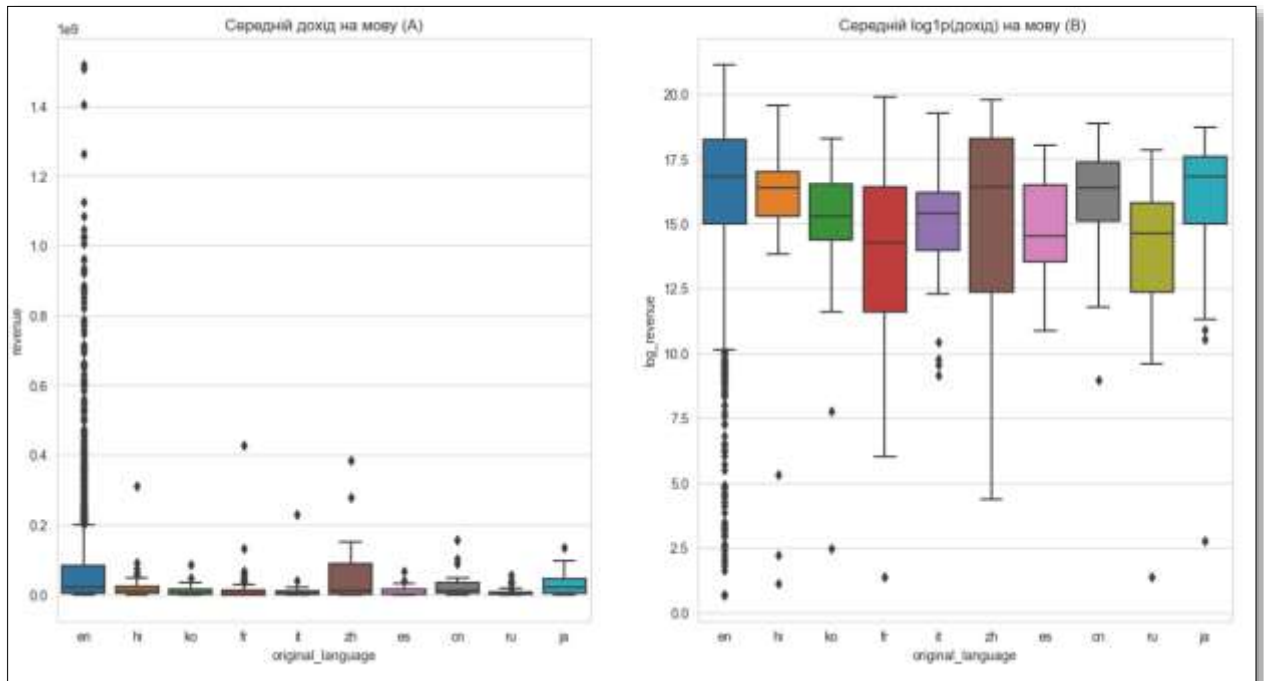


Рисунок 2.8 - Графік розподілу доходів для фільмів на різних мовах

З графіку 2.8 (А), що англійська мова має значно більший прибуток порівняно з іншими мовами. Проте після логарифмічно трансформування, рисунок 2.8, видно, що інші мови також створюють прибуток порівнянний з англійською мовою. Однак саме англійські фільми лідирують на діаграмі доходів.

#### 2.1.4 Найбільш вживані слова в фільмах.

Далі отримаємо ілюстрацію найбільш вживаних слів в назвах фільмів, за допомогою технології *WordCloud*. *WordCloud* – це метод візуалізації даних, який використовується для представлення текстових даних, де розмір слова вказує на його частоту чи важливість [12]. Для даного експерименту використано *Python* бібліотеку *wordcloud* [13]. Результат роботи *WordCloud* зображено на рисунку 2.9, з нього видно, що найбільш вживаними словами в назвах фільмів є – *Man, Last, Love, La, Life, Death*.





допоможе налагодити класифікатор машинного навчання, а також допомагає пояснити прогноз. Це допоможе виявити слова, які найбільше впливають до прибуток фільму. На рисунку 2.11 зображено вагу впливу деяких слів на прибуток фільму.

Weight <sup>2</sup>	Feature
+13.074	to
+10.131	bombing
+9.981	the
+9.777	complications
... 3858 more positive ...	
... 3315 more negative ...	
-9.281	politicians
-9.391	18
-9.481	violence
-9.628	escape and
-9.716	life they
-10.021	ones
-10.111	sally
-10.291	attracted to
-10.321	who also
-10.421	casino
-10.614	receiving
-10.759	kept
-12.139	and be
-12.939	campaign
-13.858	mike
-15.273	woman from

Рисунок 2.11 – Ваги впливу слів з описів до фільму на прибуток

З рисунка 2.11 видно, що слова можуть впливати як позитивно так і негативно на прибуток фільму. Такі слова як *to*, *bombing*, *complication* мають позитивний вплив, а такі слова як *politicians*, *18*, *violence* мають негативний вплив на прибуток.

Аналогічним способом, отримуємо ваги слів з назв фільмів. Результат зображено на рисунку 3.12.



Contribution'	Feature
+12.762	<BIAS>
+1.302	the chaos
+0.917	to
+0.874	fred
+0.760	chaos and
+0.633	s home
+0.555	return to
+0.504	home
+0.462	her job
+0.456	the
+0.390	creates
+0.355	escape from
+0.354	her
+0.321	childhood
+0.307	husband
+0.278	mother s
+0.221	from
+0.196	after
+0.196	her mother
+0.179	to win
+0.135	elizabeth
+0.129	s
+0.127	marriage
+0.108	up
+0.093	to her
+0.089	and her
+0.088	husband and
+0.074	between
+0.071	that
+0.068	of
+0.060	returns to
+0.057	and
+0.050	when
+0.047	win
+0.024	losing her
+0.003	breaks
-0.042	she
-0.057	between the
-0.062	in
-0.086	her husband
-0.100	job
-0.113	losing
-0.130	attempts
-0.145	after her
-0.232	friend
-0.255	returns
-0.261	escape
-0.284	attempts to
-0.290	mother
-0.327	to escape
-0.419	back
-0.478	job in
-0.481	from the
-0.504	return
-0.695	mayhem
-0.913	and return
-0.927	chaos

Рисунок 2.12 – Ваги впливу слів з назв фільмів на прибуток

### 2.1.6 Вплив дати виходу на прибуток кінофільму

Побудуємо діаграми з розбиттям доходу по місяцям, кварталам року, днів тижня, щоб з'ясувати чи впливає дата виходу фільму на його прибуток. На

рисунка 2.13, 2.14, 2.15 та 2.16 зображено розбиття доходу за місяцем, кварталом, днем тижня відповідно.

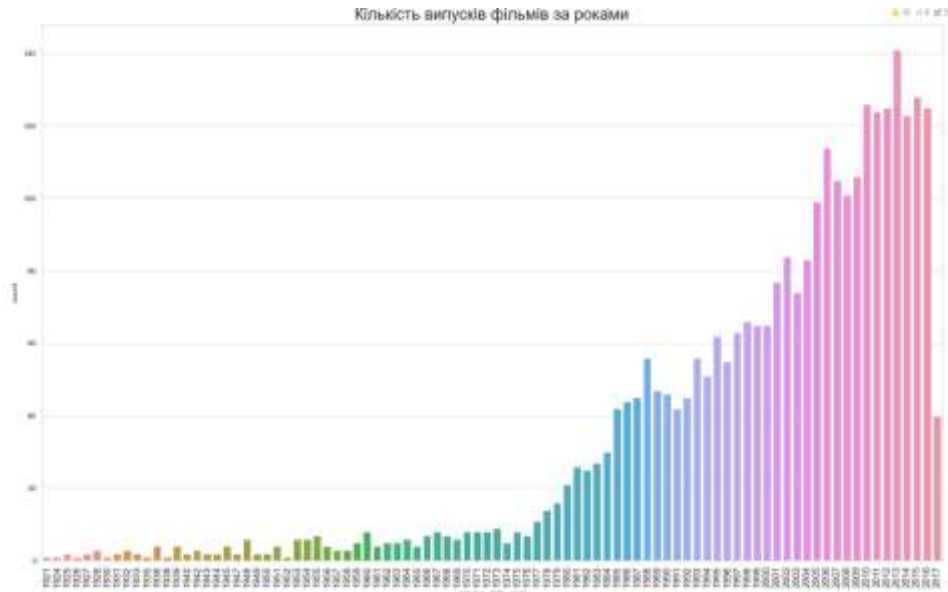


Рисунок 2.13 – Кількість випусків фільмів за роками

На рисунку 2.13 видно, що з 2000-х років спостерігається значний стрибок в кількості випущених фільмів. З діаграми випливає, що в 2013 році було випущено найбільше фільмів, тобто 140 фільмів на рік.

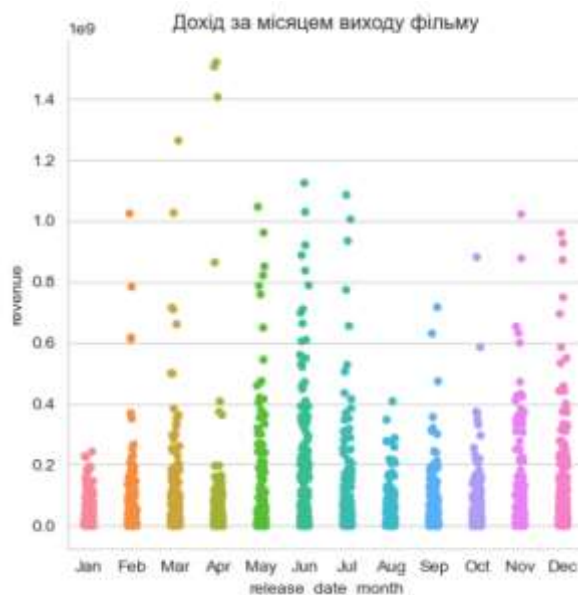


Рисунок 2.14 – Графік доходу за місяцем виходу фільму

З рисунка 2.14 випливає, що фільм, випущений у квітні, має максимальний прибуток, тоді як фільм, випущений у січні, має менший прибуток порівняно з іншими місяцями.

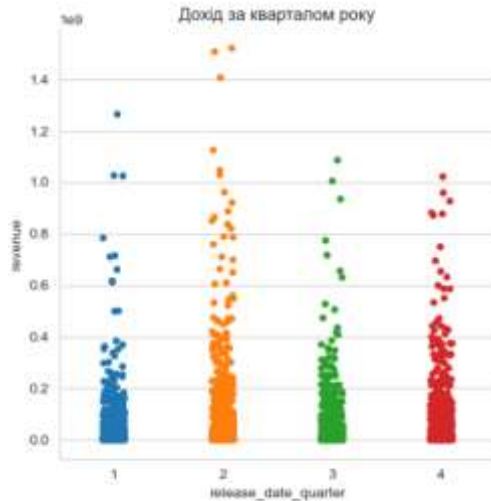


Рисунок 2.15 – Графік доходу за кварталом року

З рисунка 2.15 випливає, що фільм, випущений у другому кварталі (квітень-червень), має більший прибуток порівняно з фільмом, випущеним у іншому кварталі.

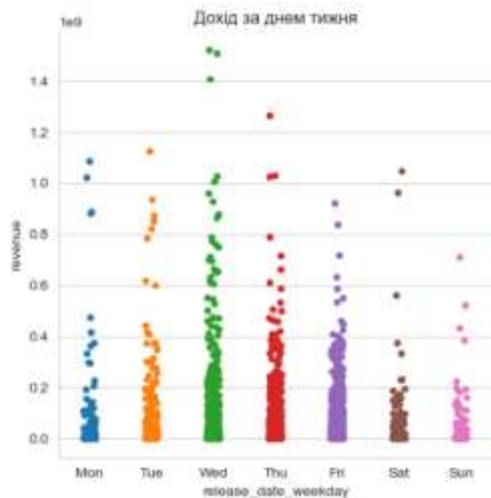


Рисунок 2.16 – Графік доходу за днем тижня

З рисунка 2.16, випливає що фільми, випущений в середу та четверг має більший прибуток. Але це може корелювати з тим, що найбільш очікувані прем'єри зазвичай виходять в четверг, а попередній показ в середу.

### 2.1.7 Співвідношення між хронометражем кінофільму та доходом.

На рисунку 2.17 зображено розподіл доходу між фільмами з різним хронометражем.

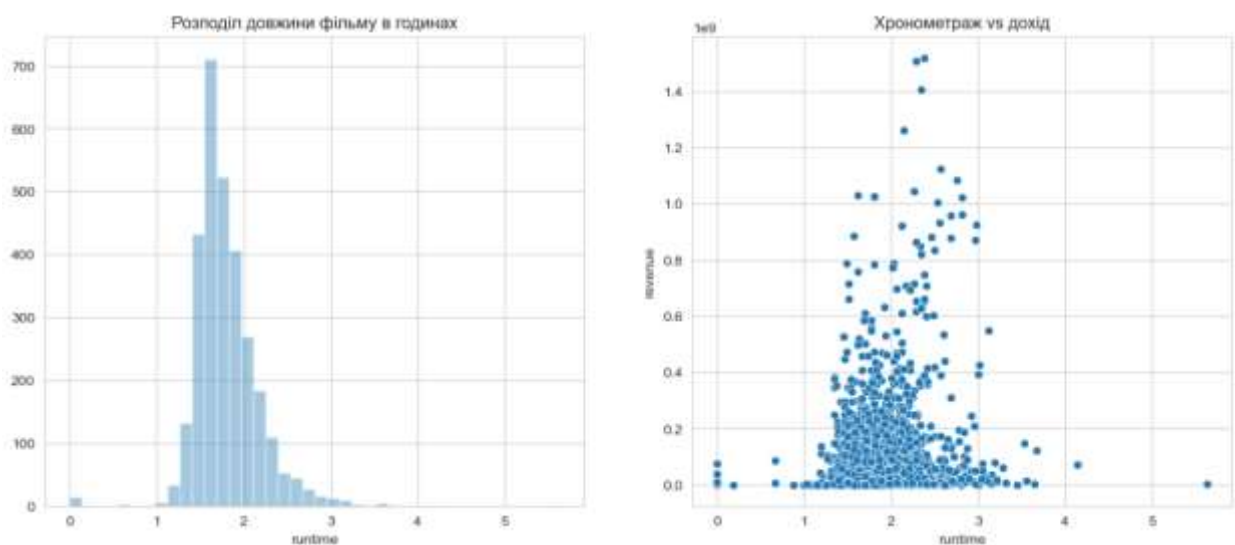


Рисунок 2.17 – Співвідношення хронометражу фільму з його доходом

З рисунка 2.17 випливає, що більшість фільмів триває від 1 до 3 годин. І фільм, який припадає на цей проміжок часу, має найвищий прибуток.

## 2.2 Розробка моделі прогнозування успішності кінофільму

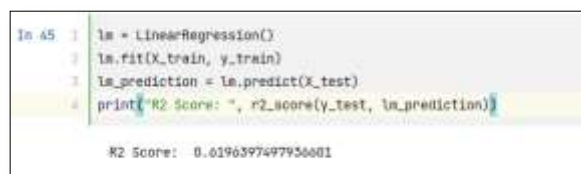
Для тренування та валідації моделі прогнозування успішності кінофільму вхідний набір даних розбитий на дві секції – дані для тренування моделі та дані для валідації. Ознаки кінофільмів, що впливають на його прибуток та зміні передбачення. Ознаки фільмів, що будуть брати участь в розробці моделі прогнозування успішності кінофільму, перелічені у таблиці 2.2.

Таблиця 2.2 – ознаки для побудови моделі прогнозування успішності фільмів

Назва ознаки	Тип даних	Опис
budget	integer	Бюджет фільму
popularity	float	Коефіцієнт популярності
runtime	integer	Хронометраж
revenue	integer	Прибуток
log1p(revenue)	float	log1p(Прибуток)
log1p(budget)	float	log1p(Бюджет)
has_homepage	boolean	Наявність домашньої сторінки
release_date_year	integer	Рік виходу
release_date_weekday	integer	День тижня виходу
release_date_month	integer	Місяць виходу
release_date_day	integer	День місяць виходу
release_date_quarter	integer	Квартал виходу

### 2.2.1 Метод лінійної регресії

Побудуємо модель використовуючи лінійну регресію. Для цього використаємо *Python* пакет *scikit-learn* та клас з цього пакету *LinearRegression*. Для перевірки точності використано коефіцієнт детермінації (*r2\_core*) [16]. На рисунку 2.18 зображено застосування лінійної регресії на ознаках виділених в таблиці 2.2.



```
In 45 | 1 | lr = LinearRegression()
      | 2 | lr.fit(X_train, y_train)
      | 3 | lr_prediction = lr.predict(X_test)
      | 4 | print("R2 Score: ", r2_score(y_test, lr_prediction))

R2 Score: 0.6196397497936681
```

Рисунок 2.18 – Застосування лінійної регресії

З рисунка 2.18 видно, що точність лінійної регресії складає **61.96%**.

### 2.2.2 Метод Random Forest.

Побудуємо модель використовуючи random forest. Для цього використаємо *Python* пакет *scikit-learn* та клас з цього пакету *RandomForestRegressor*. Для перевірки точності використано коефіцієнт детермінації. На рисунку 2.19 зображено застосування random forest на ознаках виділених в таблиці 2.2.

```
In 40 | from sklearn.ensemble import RandomForestRegressor
      |
      | RF_model = RandomForestRegressor(random_state=0, n_estimators=500, max_depth=10)
      | RF_model.fit(X_train, y_train)
      |
      | y_hat = RF_model.predict(X_test)
      | print ("R2 score:", r2_score(y_hat, y_test))
      |
      | R2 score: 0.5390510534381079
```

Рисунок 2.19 - Застосування random forest

З рисунка 2.18 видно, що коефіцієнт детермінації складає **53.96%**. Однієї з переваг *RandomForestRegressor* в *scikit-learn*, що можна отримати коефіцієнти впливу ознак на результат передбачення. На рисунку 2.20 зображено діаграму впливу ознак на результат передбачення.

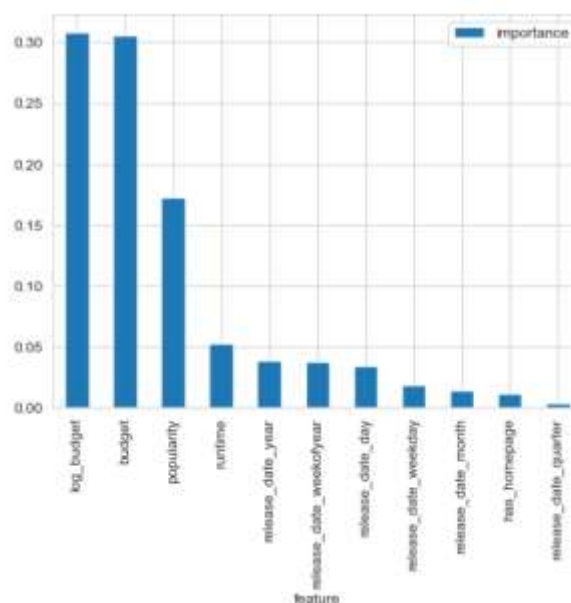


Рисунок 2.20 – Діаграма впливу ознак на результат передбачення

### 2.2.3 Метод Gradient Boosting.

Побудуємо модель використовуючи `gradient boosting`. Для цього використаємо *Python* пакет *scikit-learn* та клас з цього пакету *GradientBoostingRegressor*. Для перевірки точності використано коефіцієнт детермінації та середньоквадратичну похибку (`mse`) [17]. На рисунку 2.21 зображено застосування *GradientBoostingRegressor* на ознаках виділених в таблиці 2.2.

```

1 from sklearn import ensemble
2 params = {'n_estimators': 100, 'max_depth': 4, 'min_samples_split': 2,
3          'learning_rate': .01, 'loss': 'ls'}
4 clf = ensemble.GradientBoostingRegressor(**params)
5 predictions2 = clf.fit(X_train, y_train)
6 training_score = clf.score(X_train, y_train)
7 print("Training Score: (training_score)")
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
262
```

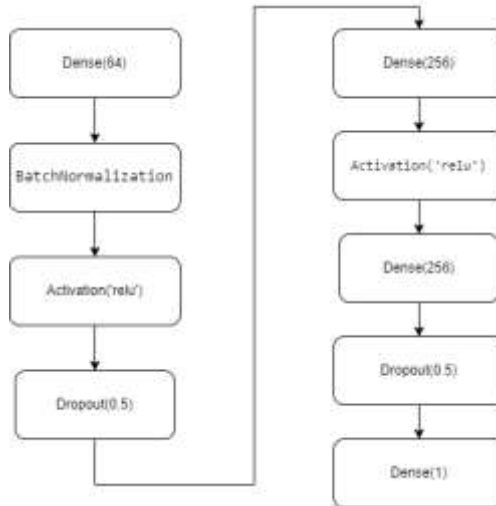


Рисунок 2.22 – структура регресійної моделі прогнозування успішності кінофільму

Розглянемо шари, присутні в нейронній мережі:

- *Dense (256)* – звичайний повно зв'язний шар, з 256 нейронами [19].
- *BatchNormalization* – застосовує перетворення, щоб підтримувати середнє значення вихідних значень шару до проміжку від 0 до 1 [20].
- *Activation('relu')* - функція активації *relu*, графік функції зображений на рисунку 2.23 [21].
- *Dropout* – випадково встановлює деякі елементів вхідної матриці рівним нулю, з вказаною частотою [22].

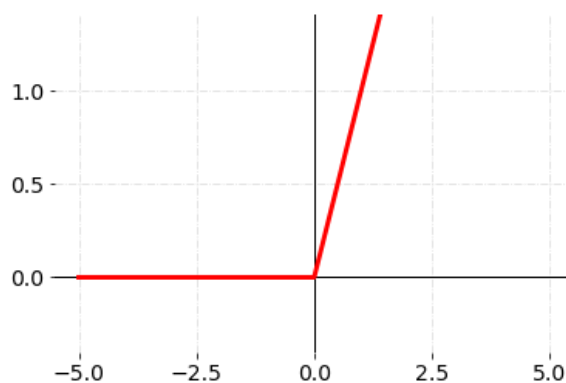


Рисунок 2.23 – Графік функції активації *relu*



Важливо зазначити, що шар *BatchNormalization* працює по різному під час навчання і в режимі прогнозування мережі. Під час навчання даний шар нормалізує дані використовуючи середнє значення та стандартне відхилення вхідних даних. Під час прогнозування шар нормалізує значення використовуючи тієї самі параметри, що і для тренування.

Нейронна мережа закінчується шаром *Dense*, що має лише один нейрон, без функції активації. Це типова конфігурація для скалярної регресії (метою якої є прогнозування одного значення на неперервній числовій прямій).

Застосування функції активації могло б обмежувати діапазон вихідних значень: наприклад, якщо в останньому шарі застосувати сигмоїдну функцію активації, мережа навчалась би прогнозувати тільки значення в проміжку між 0 та 1.

Модель компілюється з функцією втрат – середня квадратична похибка. Ця функція широко використовується в задачах регресії.

### **2.3 Висновок до розділу 2**

Проаналізовано та досліджено набір даних TMDb. Побудовано графіки кореляцій між ознаками кінофільму та його доходу. Відокремлено основні ознаки, що впливають на прибуток кінофільму.

У даному розділі було проаналізовано різні підходи для побудови моделі передбачення успішності кінофільму. Розглянуто наступні алгоритми: лінійна регресія, random forest, gradient boosting.

Запропоновано власну модель прогнозування успішності кінофільму на основі нейронної мережі.

## **3 РОЗРОБКА ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ**

### **3.1 Розробка структури програмного забезпечення прогнозування успішності кінофільму**

Для того, щоб визначитись з структурою програмного забезпечення для прогнозування успішності кінофільмам потрібно виділити вимоги до системи:

- наївність інтерфейсу користувача;
- база даних кінофільмів;
- можливість вносити дані про фільмі та редагувати їх;
- можливість імпортувати фільми з існуючих джерел;
- прогнозування успішності кінофільму повинно бути відносно швидким;
- результат прогнозування успішності кінофільму повинен бути зрозумілим користувачу.

Також система повинна виділяти бізнес логіку від інтерфейсу користувача. У такий спосіб інтерфейс може бути реалізованим різними способами, в залежності від вимог системи та бізнес плану. Наприклад інтерфейс може бути реалізованим у наступні способи:

- веб-сайт;
- мобільний додаток;
- у вигляді командного рядка;
- десктопний додаток;
- бот в месенджері;
- інтеграції з іншими сервісами.

З вище наведених вимог можна згенерувати загальну структуру системи

прогнозування успішності кінофільмів, що зображена на рисунку 3.1

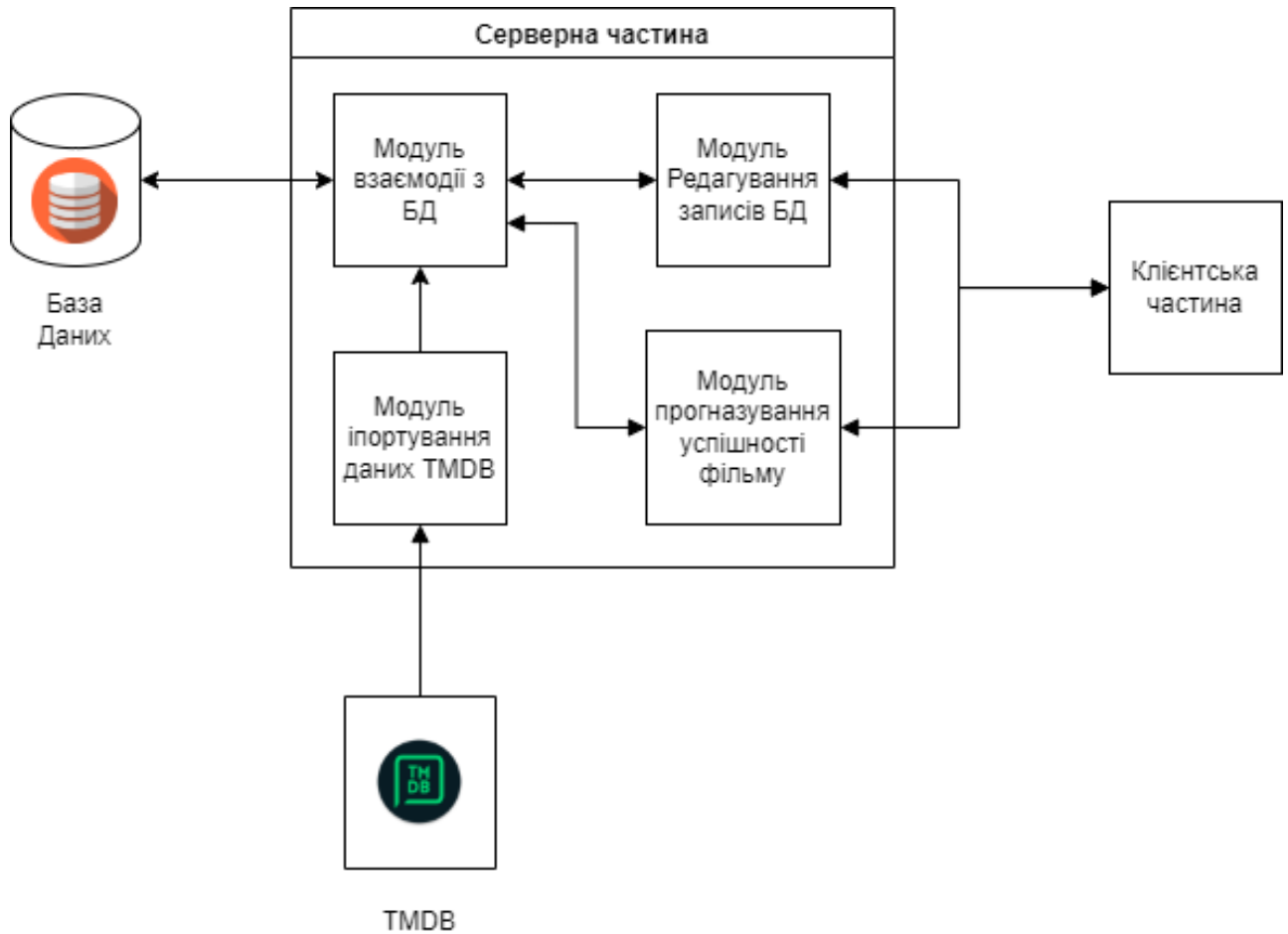


Рисунок 3.1 – Загальна структура програмного забезпечення прогнозування успішності фільмів

Джерелом імпортування даних кінофільмів, обрано найбільшу базу даних фільмів TMD. Обраний саме цей сервіс, оскільки він надає доступ до найширшої бібліотеки кінофільмів, містить різноманітні мета дані про фільмів, має інтерфейси інтеграцій.

Система прогнозування кінофільмів буде виконана у вигляді клієнт-серверного програми. Таким чином, логіка зберігання фільмів та прогнозування буде інкапсульована в серверній частині, а логіка інтерфейсу користувача в клієнтській частині системи.

Для реалізації серверної частини потрібно обрати архітектуру. Для налаштування спілкування між двома програма розробники будують мости –

програмний інтерфейс програма (API), щоб дозволити одній системі отримати доступ до даних або функціональності іншої.

Для швидкої та масштабної інтеграції програм *API* реалізується з використанням протоколів та/або специфікації для визначення семантики та синтаксису повідомлень, що передаються по зв'язку. Ці специфікації складають архітектуру *API*. З часом були створенні різні архітектури стилі *API*. Кожен з них має свої стандарти обміну даними. У таблиці 3.1 наведено порівняльну характеристику існуючих архітектур серверних додатків [23].

Таблиця 3.1 порівняльна характеристика архітектури серверних додатків

	RPC	SOAP	REST	GraphQL
Опис	Виклик локальних процедур	Обгорнута структура повідомлень	Дотримання 6 архітектурних обмежень	Схема та система типів
Формат	JSON, XML, Protobuf, Thrift, FlatBuffres	XML	XML, JSON, HTML, TEXT	JSON
Складність	Проста	Складна	Проста	Середня

Оскільки *REST* архітектура відносно проста в реалізації там підтримує багато форматів спілкування, було обрано саме цю архітектуру клієнт-серверної частини.

### 3.2 Проектування бази-даних інформаційної технології прогнозування успішності фільмів

Для зберігання інформації, що буде використовуватись для прогнозування успішності кінофільму, потрібні наступні сутності:

- фільм – movie;

- жанр – genre;
- студія – production\_company;
- країна виробник – production\_country;
- ключове слово – key\_word;
- актор – actor;
- знімальний персонал – credit.

**Знімальний персонал** – сутність, що відображає людей задіяних в зйомці фільму, таки-як: режисер, сценарист монтажер і тд. В таблиці 3.1 відображено поля сутності – знімальний персонал.

Таблиця 3.2 – Поля сутності знімального персоналу (credit)

Назва поля	Тип даних	Опис
credit_id	integer, primary key	Індивікатор
department	string	Тип департаменту: Directing, Production ...
gender	enum(male, female)	Стать
job	string	Опис посади
name	string	Ім'я

**Жанр** – сутність, що відображає жанр фільму. В таблиці 3.2 відображено поля сутності – жанр.

Таблиця 3.3 – Поля сутності жанр (genre)

Назва поля	Тип даних	Опис
genre_id	integer, primary key	Індивікатор
name	string	Ім'я жанра

**Студія** – сутність, що студія розробки фільму. В таблиці 3.3 відображено

поля сутності – **студія**.

Таблиця 3.4 – Поля сутності студія (production\_company)

Назва поля	Тип даних	Опис
id	integer, primary key	Індивікатор
name	string	Ім'я студії

**Країна виробник**– сутність, що студія розробки фільму. В таблиці 3.4 відображено поля сутності – країна виробник.

Таблиця 3.5 – Поля сутності студія (production\_country)

Назва поля	Тип даних	Опис
Id	integer, primary key	Індивікатор
Name	string	Назва країни

**Ключове слово** – ключове слово, що використовується для характеристики кінофільму. В таблиці 3.5 відображено поля сутності – країна виробник.

Таблиця 3.6 – Поля сутності ключове слово (key\_word)

Назва поля	Тип даних	Опис
Id	integer, primary key	Індивікатор
Name	string	Назва слова

**Актор** – актор, що приймає участь у фільмі. В таблиці 3.6 відображено поля сутності – актор.

Таблиця 3.7 – Поля сутності актор (cast)

Назва поля	Тип даних	Опис
id	integer, primary key	Індивікатор
Назва поля	Тип даних	Опис
character	string	Ім'я персонажа
gender	Enum(male ,female)	Стать
name	Enum(male ,female)	Ім'я актора

**Фільм.** В таблиці 3.7 відображено поля сутності – фільм.

Таблиця 3.8 – Поля сутності фільм (movie)

Назва поля	Тип даних	Опис
id	integer, primary key	Індивікатор
budget	integer	Бюджет
genre_ids	List[integer], зовнішній ключ на genre	Список жанрів
homepage	string	Посилання на головну сторінку
Imdb_id	string	Індивікатор IMDB
original_language	string	Мова фільму
Назва поля	Тип даних	Опис
original_title	string	Заголовок фільму
popularity	float	Оцінка популярності
production_companies	List[integer], зовнішній ключ на product_company	Компанії виробники
production_countries	List[integer], зовнішній ключ на product_country	Країни виробники

Продовження таблиці 3.9.

release_date	Datetime	Дата релізу
runtime	Integer	Хронометраж фільму
status	String	Статус фільму
keywords	List[integer], зовнішній ключ key_wrod	Список ключових слів
Cast	List[integer], зовнішній ключ actor	Список акторів
Crew	List[integer], зовнішній ключ credit	Персонал
revenue	Integer	Прибуток

На основі вище наведених сутностей, можна згенерувати загальну схему зв'язків бази даних, що зображено на рисунку 3.2.

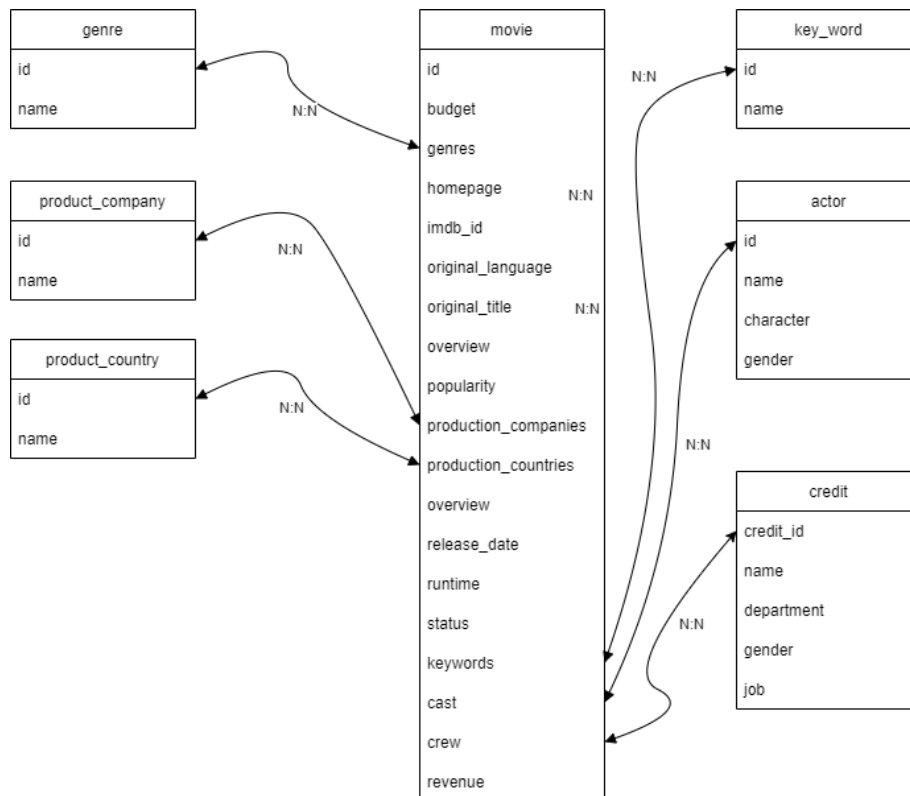


Рисунок 3.2 – Загальна схема бази-даних кінофільмів

### 3.3 Розробка алгоритму модуля взаємодії з базою даних кінофільмів



Модуль взаємодії з базою даних кінофільмів відповідає за програмне з'єднання з базою даних. Модуль повинен виконувати наступний функціонал:

- додавання нового фільму
- редагування існуючого фільму
- видалення фільму
- отримання фільмів

На рисунку 3.3 зображено схему алгоритму роботи модуля взаємодії з базою даних.

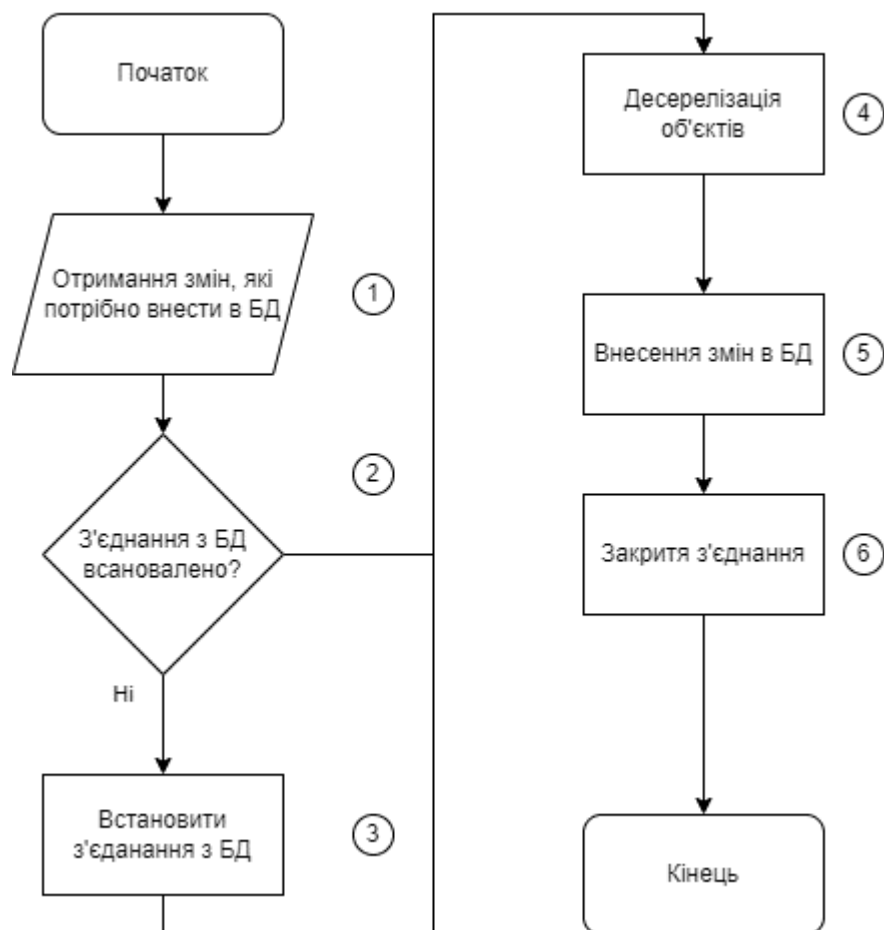


Рисунок 3.3 – Схема алгоритму модуля взаємодії з базою даних

Опишемо даний алгоритм:

1. Отримання змін, які потрібно ввести в базу даних.
2. Перевірка, чи встановлене з'єднання з базою даних.
3. Встановлення з'єднання з базою даних.
4. Десерілізація об'єктів – перетворення типів даних об'єктів мови програмування в типи даних бази даних.
5. Закриття з'єднання з базою даних.

### 3.4 Розробка алгоритму модуля імпортування кінофільмів з TMDB

Основною задачею модуля імпортування фільмів – є імпортування фільмів з найбільшої бази даних фільмів TMDB та в подальшому зберігання отриманих даних в локальній базі. При цьому модуль повинен вміти об'єднувати існуючі дані про фільми з даними отриманими з TMDB. На рисунку 3.4 зображено схему алгоритму модуля.

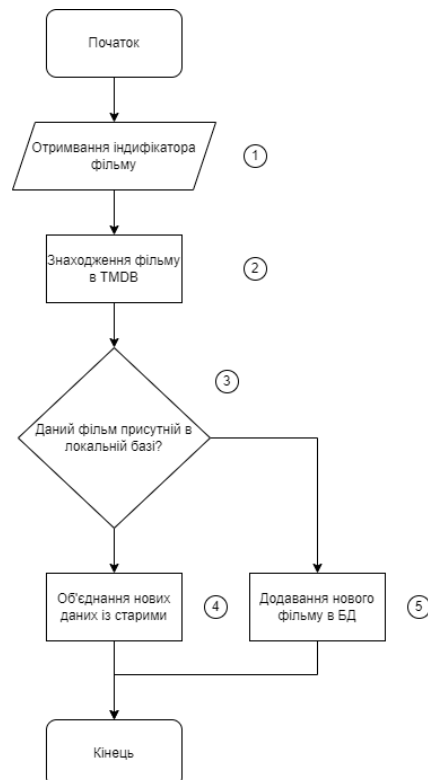


Рисунок 3.4 – Схема алгоритму модуля імпортування фільмів з TMDB

Опишемо даний алгоритм:

1. Отримання індикатора фільму.
2. Знаходження фільму в базі даних TMDb.
3. Перевірка чи в системі уже є даний фільм.
4. Об'єднання старої інформації про фільм з новою.
5. Додавання нового фільму в базу даних.

### 3.5 Розробка алгоритму модуля редагування записів бази даних фільмів

Основною задачею даного модуля – надавання інтерфейсу редагування записів бази даних фільмів для клієнтської частини. Оскільки в розділі 3.1 було обрано REST архітектуру для серверної частини, клієнт буде взаємодіє з сервером через HTTP запити. Для цього сервер повинен підтримувати усі основні HTTP методи для кожної сутності бази даних, що були описані в розділі 3.2. В таблиці 3.2 перелічено основні HTTP запити, які повинен підтримувати сервер, для надання клієнту можливості взаємодіяти з записами в базі даних. Усі вхідні та вихідні об'єкти будуть приставлені у вигляді JSON об'єктів.

Таблиця 3.10 – Основні HTTP запити модуля редагування записів БД

Назва	Метод	Шлях	Вхідні дані	Вихідні дані
Отримання сутності по індикатору	GET	/:id	Індифікатор сутності	Об'єкт сутності
Отримання усіх сутностей	GET	/	Фільтр	Об'єкти сутностей

Продовження таблиці 3.10

Назва	Метод	Шлях	Вхідні дані	Вихідні дані
Додавання сутності	POST	/	Об'єкт сутності	Об'єкт сутності
Оновлення усієї сутності	PUT	/:id	Індифікатор та об'єкт сутності	Об'єкт сутності
Оновлення окремих полів сутності	PATCH	/:id	Індифікатор та об'єкт сутності	Об'єкт сутності
Видалення сутності	DELETE	/:id	Індифікатор сутності	

### 3.6 Розробка алгоритму модуля прогнозування успішності кінофільму

Модуль прогнозування успішності кінофільму безпосередньо відповідає за надання прогнозу доходу фільму на основі ознак фільму. На рисунку 3.5 зображено схему алгоритму прогнозування фільму.



Рисунок 3.5 – Схема алгоритму модуля прогнозування успішності кінофільму

Опишемо алгоритм:

1. Отримання даних фільму.
2. Виділення ознак, що були визначені в розділі 2.
3. Нормалізація числових даних.
4. Тонізація текстових даних, так як модель прогнозування може працювати тільки чи чистовий значеннями.
5. Завантаження моделі прогнозування з розділу 2.
6. Застосування моделі для отримання прогнозу.
7. Інтерпретація прогнозу.

### **3.7 Висновок до розділу 3**

У розділі наведено основні кроки роботи інформаційної технології прогнозування успішності кінофільму. Обрано клієнт-серверну архітектуру для програмного забезпечення прогнозування успішності кінофільму. Розбито програмний засіб на модулі, а саме: модуль взаємодії з базою даних, модуль імпортування кінофільмів, модуль редагування записів в базі даних, модуль прогнозування успішності кінофільму. Описано роботу та алгоритм кожного модуля. Спроектовано базу даних для програмного забезпечення прогнозування успішності кінофільмів.

## 4 ПРОГРАМНА РЕАЛІЗАЦІЯ ІНФОРМАЦІЙНОЇ ТЕХНОЛОГІЇ ПРОГНОЗУВАННЯ УСПІШНОСТІ КІНОФІЛЬМУ

### 4.1 Обґрунтування вибору мови програмування

Серверна частина буде розроблена на мові програмування *Python 3.10*, використовуючи фреймворк *FastAPI*. *Python* було обрано так як він надає широкий спектр потужних інструментів та бібліотек для розробки серверної частини додатка, обробки зображень, роботи з великим обсягом даних, навчання, роботи з багатовимірними матрицями, тренування та використання моделей машинного навчання. *FastAPI* – це сучасний, швидкий веб-фреймворк для створення *API*-інтерфейсів з *Python 3.6+* на основі стандартних підказок типів *Python* [24]. Основні особливості *FastAPI*:

Висока продуктивність, нарівні з *NodeJS* та *Go*;

Збільшення швидкості розробки приблизно на 200% - 300 %;

Зменшення кількості помилок приблизно на 40 %, спричиненою розробником;

Розробка та навчання моделей машинного навчання відбувалася з використанням бібліотек *TensorFlow* та *Keras* - високопродуктивні бібліотеки для написання моделей машинного навчання. *TensorFlow* - відкрита програмна бібліотека для машинного навчання цілій низці задач, розроблена компанією *Google* для задоволення її потреб у системах, здатних будувати та тренувати нейронні мережі для виявлення та розшифрування образів та кореляцій, аналогічно до навчання й розуміння, які застосовують люди. *Keras* надає зручний та мінімалістичний інтерфейс над *TensorFlow* та *TensorFlow 2.0*, що пришвидшує написання час моделей та полегшує процес дослідження моделей машинного навчання.

Клієнтська частина буде розроблена на мові програмування *Dart*, використовуючи фреймворк *Flutter*. *Flutter* було обрано, оскільки він надає

програмісту можливість створювати красиві та нативні програми, задіяючи мінімальні зусилля. За допомогою *Flutter* можна створювати кросплатформені програми, використовуючи одну кодову базу [25]. *Flutter* підтримує наступні платформи:

- Web;
- Windows;
- Linux;
- MacOS;
- iOS;
- Android.

#### **4.2 Обґрунтування вибору середовища програмування**

Для розробки серверної частини інформаційної технології проказування успішності кінофільму, було обрано середовище «*PyCharm*». *PyCharm* забезпечує інтелектуальне завершення та підказки при написанні коду, перевірку коду, оперативне виділення помилок та підказки для виправлення помилок, а також автоматичний рефакторинг коду і широкі можливості навігації. Колекція інструментів *PyCharm* з коробки включає вбудований відладчик і тестовий прогін; Профілювальник Python; вбудований термінал; інтеграція з основними засобами контролю версій і вбудованими інструментами для бази даних; можливості віддаленої розробки з віддаленими інтерпретаторами; вбудований ssh-термінал; і інтеграція з *Docker* і *Vagrant*. *PyCharm* інтегрується з *IPython Notebook*, має інтерактивну консоль *Python* і підтримує *Anaconda*, а також кілька наукових пакетів, включаючи *Matplotlib* і *NumPy* [26]. Через вище вказані характеристики *PyCharm* є одним із провідних середовищ програмування на мові *Python*. Підтримка математичних модулів та *Anaconda* – безкоштовний дистрибутив пакетів для машинного навчання,

забезпечує комфортну розробку моделей для машинного навчання. На рисунку 4.1 зображено інтерактивне середовище розробки для мови програмування *Python* – *PyCharm*.

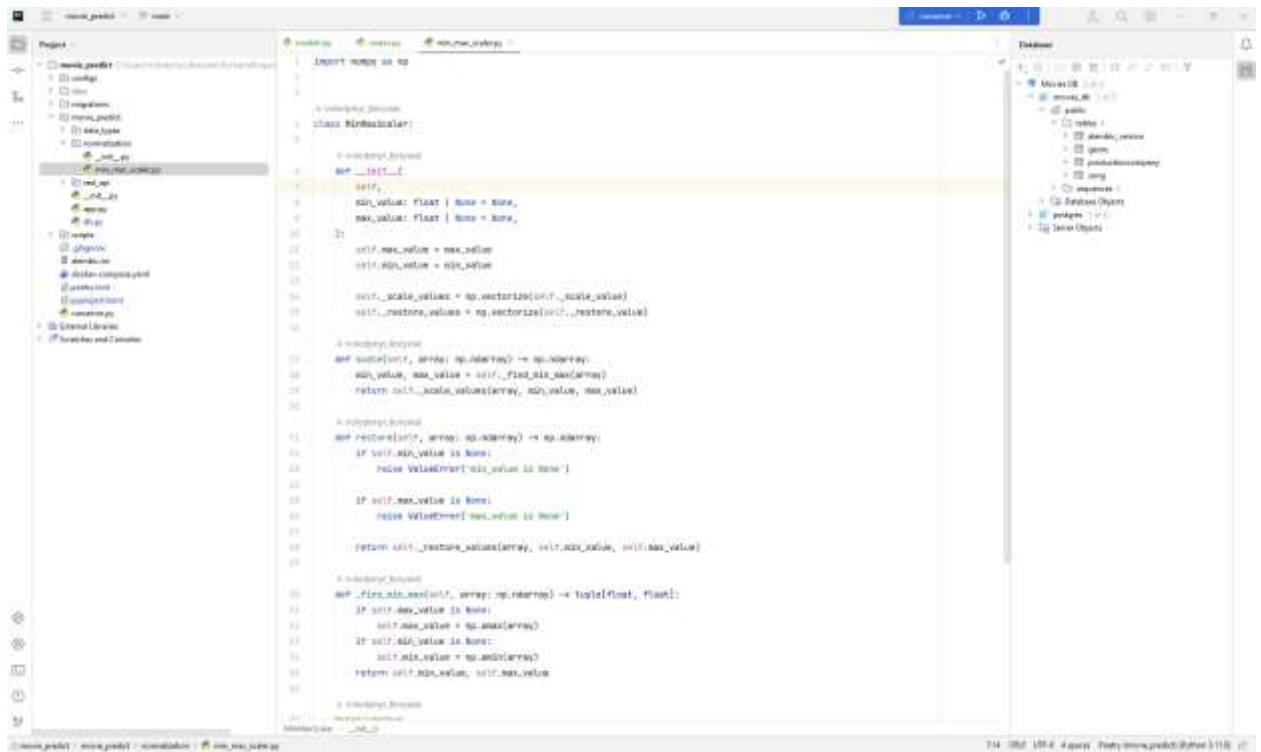


Рисунок 4.1 – середовище програмування *PyCharm*

Для розробки серверної частини інформаційної технології проказування успішності кінофільму, було обрано середовище «*IntelliJ Idea*» [27]. *IntelliJ Idea* має такі самі переваги, що і *PyCharm*, окрім підтримки математичних пакетів. Проте має підтримку мови програмування *Dart*, та фреймворку *Flutter*. Надає зручний інтерфейс для запуску емуляторів пристроїв, що надає можливість писати та тестувати програми на різних платформах, такі як:

- Windows.
- Linux.
- Android.
- iOS.



На рисунку 4.2 зображено інтерактивне середовище розробки – *IntelliJ Idea*.

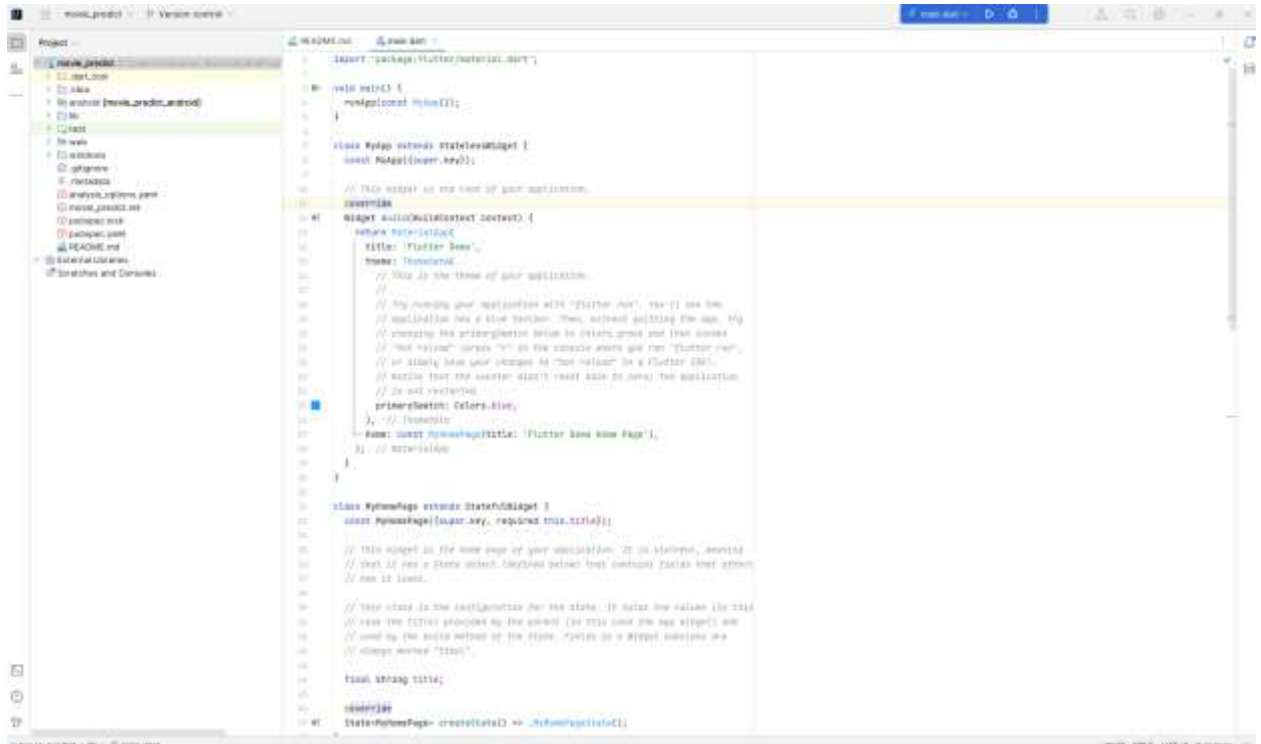


Рисунок 4.2 – середовище програмування *IntelliJ Idea*.

### 4.3 Програмна реалізація модулів інформаційної технології прогнозування успішності кінофільмів

#### 4.3.1 Програмна реалізація модуля взаємодію з базою даних

Модуль взаємодії з базою даних реалізований на основі Python бібліотеки `sqlmodel`. `Sqlmodel` – надає зручний інтерфейс опису сутностей за допомогою класів та автоматичне створення SQL запитів на основі моделей. Інтерфейс `Repo` описує взаємодію з базою даних для кожної сутності. На рисунку 4.3 зображено фрагмент коду `Repo`.

```

class Repo(Generic[_Model], ABC):

    def __init__(self, session: Session):
        self._session = session

    @abstractmethod
    def get(self, pk: int) → _Model:
        pass

    @abstractmethod
    def get_many(self, *where) → Iterable[_Model]:
        pass

    @abstractmethod
    def delete(self, pk: int) → NoReturn:
        pass

```

Рисунок 4.3 – Фрагмент коду інтерфейсу Repo

На рисунку 4.4 зображено фрагмент реалізації інтерфейсу Repo для сутності Movie.

```

class MovieRepo(Repo[Movie]):

    def get(self, pk: int) → Movie:
        query = select(Movie).where(Movie.id == pk)
        movie = self._session.execute(query).first()
        return movie

    def get_many(self, *where) → Iterable[Movie]:
        query = select(Movie).where(*where)
        return self._session.execute(query)

    def delete(self, pk: int) → NoReturn:
        query = delete(Movie).where(Movie.id == pk)
        self._session.execute(query)

```

Рисунок 4.4 – Фрагмент коду класу MovieRepo

На рисунку 4.5 зображено реалізацію сутності Movie за допомогою sqlalchemy.

```

class Movie(SQLModel, table=True):
    id: int = Field(primary_key=True)
    budget: int
    genre_ids: list[ForeignKey('genre.id')] = []
    homepage: str | None
    imdb_id: str
    original_language: str
    original_title: str
    overview: str
    popularity: float
    production_companies: list[ForeignKey('production_company.id')] = []
    production_countries: list[ForeignKey('production_country.id')] = []
    release_date: dt.datetime
    runtime: int
    keywords: list[ForeignKey('keyword.id')] = []
    cast: list[ForeignKey('actor.id')] = []
    crew: list[ForeignKey('credit.id')] = []
    revenue: int | None

```

Рисунок 4.5 – Реалізація сутності Movie за допомогою sqlalchemy.

#### 4.3.2 Програмна реалізація модуля редагування записів в базі даних

Модуль редагування записів в базі даних реалізовано на основі фреймворку FastAPI, який надає можливості створювати обробники HTTP записів за допомогою декораторів. На рисунку 4.6 зображено фрагмент коду, що описує обробники HTTP запитів для редагування фільмів в базі даних.

```

@app.router.put('/', response_model=MovieInDB)
async def create_movie(movie: CreateMovie) → MovieInDB:
    movie_db = MovieInDB.parse_obj(movie)
    return await movie_db.insert()

# volodymyrb
@app.router.get('/{movie_id}', response_model=MovieInDB)
async def get_movie(movie_id: PydanticObjectId) → MovieInDB:
    return await MovieInDB.get(movie_id)

# volodymyrb
@app.router.get('/', response_model=list[MovieInDB])
async def get_movies() → list[MovieInDB]:
    return [movie async for movie in MovieInDB.find_all()]

# volodymyrb
@app.router.patch('/{movie_id}', response_model=MovieInDB)
async def update_movie(movie_id: PydanticObjectId, movie: UpdateMovie) → MovieInDB:
    movie_db = await MovieInDB.get(movie_id)
    await movie_db.set(movie.dict(exclude_none=True))

```

Рисунок 4.6 – Фрагмент коду для обробки запитів редагування фільмів.

### 4.3.3 Програмна реалізація модуля імпортування кінофільмів з *TMDB*

Модуль призначений для імпортування даних про фільми з сервісу *TMDB*, для того щоб не заповнювати усі інформація власноруч.

Модуль імпортування фільмів реалізований у вигляді *HTTP* клієнта. Клієнт створює запити на сервіс *TMDB*, вказуючи Індифікатор фільму, що потрібно знайти. Далі модуль завантажує отриману інформація в базу даних. На рисунку 4.7 зображено фрагмент коду *HTTP* клієнта, а на рисунку 4.8 фрагмент коду завантаження результату в базу даних.

```
class TmdbClient(AsyncConnection[httpx.AsyncClient]):

    _BASE_URL = Url(
        host=CONFIG.TMDB.HOST,
        username=CONFIG.TMDB.TOKEN_ID,
        password=CONFIG.TMDB.TOKEN_VALUE,
    )

    _MOVE_URL = _BASE_URL.join_path(Path('movie'))

new *
async def close(self) → bool:
    await self._connected_conn.aclose()
    return True

new *
async def _create_conn(self) → httpx.AsyncClient:
    return httpx.AsyncClient()

new *
def get_movie(self, movie_pk: str) → Movie:
    movie_obj = self._connected_conn.get(
        url=self._MOVE_URL.update_query({
            'pk': movie_pk,
        }).build(),
    )

    if not movie_obj:
        raise ValueError(f'Movie {movie_obj} is not found')

    movie = Movie.parse_obj(movie_obj)
    return movie
```

Рисунок 4.7 – Фрагмент коду модуля імпортування кінофільмів з *TMDB*

```

new *
async def main():
    unique_production_companies = {
        load_production_companies(TRAIN_DATA_PATH)
        | load_production_companies(TEST_DATA_PATH)
    }
    unique_production_companies.add('other')

    await db.init_db()
    async with db.get_session_ctx() as session:
        for company_name in unique_production_companies:
            company = ProductionCompany(name=company_name)
            session.add(company)
        await session.commit()

if __name__ == '__main__':
    asyncio.run(main())

```

Рисунок 4.8 – Фрагмент модуля завантаження даних з *TMDB* в базу даних

#### 4.3.4 Програмна реалізація модуля прогнозування успішності кінофільму

Безпосереднє прогнозування доходу фільму реалізовано в класі *Prediction*, який завантажує натреновану нейронну мережу за проганяє через неї вхідні дані кінофільму, на виході отримуємо прогнозований прибуток кінофільму. На рисунку 4.9 зображено фрагмент класу *Prediction*.

```

from decimal import Decimal

import keras
from keras import models

from configs import CONFIG
from movie_predict.rest_api.movies.models import Movie
from movie_predict.normalization.movie_normalization import prepare_movie

class Prediction:

    def __init__(self):
        self._model = self._load_model()

    def predict_revenue(self, movie: Movie) -> Decimal:
        tensor = prepare_movie(movie)
        prediction = self._model(tensor)
        return prediction[0]

    @staticmethod
    def _load_model() -> keras.Sequential:
        return models.load_model(str(CONFIG.MOVIE_PREDICTION.MODEL_PATH))

```

Рисунок 4.9 – Фрагмент модуля передбачення успішності кінофільму

Взаємодія з клієнтом реалізовано за допомогою обробника HTTP запитів - `get_predicted_revenue`. Даний обробник отримує на вхід ідентифікатор фільму, знаходить потрібний фільм в базі даних, застосовує клас `Prediction`, та повертає прогнозований прибуток кінофільму. На рисунку 4.10 зображено фрагмент реалізації обробника.

```
router = APIRouter(
    prefix='/prediction',
    tags=['Movie prediction'],
)

new *
@router.get('/revenue')
def get_predicted_revenue(
    movie_pk: int,
    session: AsyncSession = Depends(get_session),
) → dict:

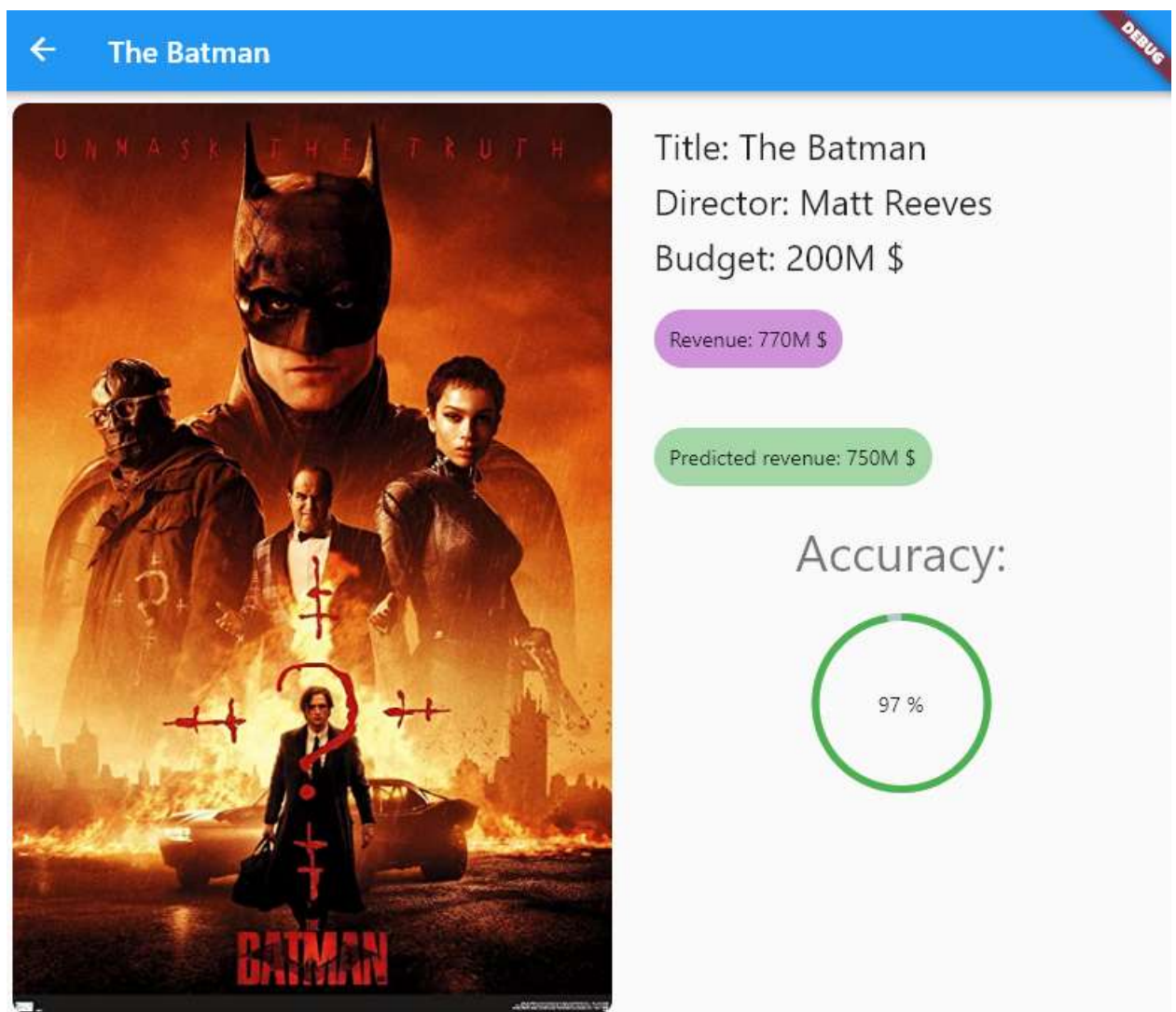
    statement = select(Movie).where(Movie.id == movie_pk)
    response = await session.execute(statement)
    movie = response.scalar_one()

    prediction = Prediction()
    return {
        'result': prediction.predict_revenue(movie),
    }
```

Рисунок 4.10 – Фрагмент коду обробника HTTP запитів, для прогнозування доходу кінофільму

#### 4.4 Тестування та аналіз результатів роботи програми прогнозування успішності кінофільму

Спочатку протестуємо в цілому роботу програми прогнозування успішності кінофільму. Вікно з результатом роботи програми зображено на рисунку 4.11.



Рисунк 4.11 – Вікно програми з результатом прогнозування успішності кінофільму

З рисунка 4.5 видно, що програма для кінофільму “The Batman” з реальним доходом в 770 мільйонів доларів США, надала прогноз доходу в 750

мільйонів. Для даного випадку, точність прогнозу досягаю 97%, що цілком відповідає поставленим вимогам.

Для доведення факту досягнення поставленої в роботі мети – підвищення точності прогнозування успішності кінофільму – було проведено тестування роботи розробленої програми та аналогічних методів. Програми-аналоги не були розглянуті так як таких немає в публічному доступі. Тестова вибірка складає 600 кінофільмів з набору даних *TMDB*. У якості метрики, для порівняння точності різних методів прогнозування, обрано середню абсолютну похибку (*mae*) [\*].

$$mae = \frac{\sum_{i=1}^n |y_i - x_i|}{n}, \quad (4.1)$$

де  $y_i$  – передбачення;  $x_i$  – справжнє значення;  $n$  – загальна кількість тестів.

Результати тестування подані в таблиці 4.1.

Таблиця 4.1 – Результати тестування розробленої програми та аналогів

Метод прогнозування	Кількість кінофільмів у тестовій вибірці	Кількість прогонів	Точність прогнозування (на основі <i>mae</i> ), %
Лінійна регресія	1000	4	61.96
Random forest	1000	4	54
Gradient Boosting	1000	4	85.37
Регресія на основі нейронної мережі (запропонований метод)	1000	4	94

Із таблиці 4.1 помітно, що розроблена програма характеризується вищою точністю проказування (94%), ніж аналоги (85%), а це означає, що



достовірність прогнозування успішності кінофільму у розробленій програмі підвищена на 9%, що свідчить про те, що мета роботи досягнута.

На рисунку 4.12 зображено фрагмент коду для обрахунку точності нейромережевого методу.

```
def validate():
    x_train, x_test, y_train, y_test, _, _, _ = samples

    k = 4
    num_val_samples = len(x_train) // k
    num_epochs = 100
    all_scores = []

    model: models.Sequential = models.load_model('trained_model.h5')
    for i in range(k):
        print(f'Processing # {i}')

        val_data = x_train[
            i * num_val_samples: (i + 1) * num_val_samples
        ]
        val_targets = y_train[
            i * num_val_samples: (i + 1) * num_val_samples
        ]

        val_mse, val_mae = model.evaluate(val_data, val_targets, verbose=0)
        all_scores.append(val_mae)

    print(all_scores)
    print(np.mean(all_scores))
```

Рисунок 4.12 – Фрагмент коду тестування нейромережевої моделі

Таким чином, після проведення порівняння розробленої програми прогнозування успішності кінофільму на основі нейромережі з аналогами, можна зробити висновок, що розроблена програма має більшу точність на 9%. Деякий ілюстративний матеріал до програми (у т.ч. скріншоти) подано в додатку В. Інструкцію користування розробленою програмою наведено у додатку Г.

#### 4.5 Висновок до розділу 4

У розділі обґрунтовано вибір мови програмування *Python* для серверної частини, та *Dart* – для клієнтської частини. У результаті було розроблено програмне забезпечення прогнозування успішності кінофільму, створену мовою програмування *Python* із застосуванням безкоштовної бібліотек: *keras*, *tensorflow*, *fastapi*, *sqlmodel*, *urlx*, *iter-model*. Було проведено тестування програми нейромережевого проказування успішності кінофільму. Аналіз результатів тестування показує, що точність прогнозування складає 94%, у той час коли точність аналогів складає 85% при навчанні на однаковій навчальній множині у кількість 2400 кінофільмів. Тобто мета магістерської кваліфікаційної роботи досягнута – точність прогнозування успішності кінофільму підвищена на 9%.

## 5 ЕКОНОМІЧНА ЧАСТИНА

Науково-технічна розробка має право на існування та впровадження, якщо вона відповідає вимогам часу, як в напрямку науково-технічного прогресу та і в плані економіки. Тому для науково-дослідної роботи необхідно оцінювати економічну ефективність результатів виконаної роботи.

Магістерська кваліфікаційна робота з розробки та дослідження «Інформаційна технологія прогнозування успішності кінофільму» відноситься до науково-технічних робіт, які орієнтовані на виведення на ринок (або рішення про виведення науково-технічної розробки на ринок може бути прийнято у процесі проведення самої роботи), тобто коли відбувається так звана комерціалізація науково-технічної розробки. Цей напрямок є пріоритетним, оскільки результатами розробки можуть користуватися інші споживачі, отримуючи при цьому певний економічний ефект. Але для цього потрібно знайти потенційного інвестора, який би взявся за реалізацію цього проекту і переконати його в економічній доцільності такого кроку.

Для наведеного випадку нами мають бути виконані такі етапи робіт:

- 1) проведено комерційний аудит науково-технічної розробки, тобто встановлення її науково-технічного рівня та комерційного потенціалу;
- 2) розраховано витрати на здійснення науково-технічної розробки;
- 3) розрахована економічна ефективність науково-технічної розробки у випадку її впровадження і комерціалізації потенційним інвестором і проведено обґрунтування економічної доцільності комерціалізації потенційним інвестором.

### **5.1 Проведення комерційного та технологічного аудиту науково-технічної розробки**

Метою проведення комерційного і технологічного аудиту дослідження за темою «Інформаційна технологія прогнозування успішності кінофільму» є

оцінювання науково-технічного рівня та рівня комерційного потенціалу розробки, створеної в результаті науково-технічної діяльності.

З появою Інтернету сильно зросла кількість інформації, з якої люди щодня стикаються. Це означає, що люди повинні орієнтуватися серед надзвичайно великої кількості доступних альтернатив, коли хочуть щось знайти. Наприклад, від вибору нового мобільного телефону або плеєра до пошуку кінофільму для вечірнього перегляду. З іншого боку виступають власники інтернет-магазинів і сервісів: вони зацікавлені в персональній рекламі і рекомендаціях кожному конкретному користувачеві, тому що такий підхід може істотно збільшити прибуток компаній. Як результат, в останні роки інтерес до розробки та вдосконалення існуючих рекомендаційних систем значно виріс.

Оцінювання науково-технічного рівня розробки та її комерційного потенціалу рекомендується здійснювати із застосуванням 5-ти бальної системи оцінювання за 12-ма критеріями, наведеними в табл. 4.1 [28].

Таблиця 5.1 – Рекомендовані критерії оцінювання науково-технічного рівня і комерційного потенціалу розробки та бальна оцінка

Бали (за 5-ти бальною шкалою)					
	0	1	2	3	4
Технічна здійсненність концепції					
1	Достовірність концепції не підтверджена	Концепція підтверджена експертними висновками	Концепція підтверджена розрахунками	Концепція перевірена на практиці	Перевірено працездатність продукту в реальних умовах
Ринкові переваги (недоліки)					
2	Багато аналогів на малому ринку	Мало аналогів на малому ринку	Кілька аналогів на великому ринку	Один аналог на великому ринку	Продукт не має аналогів на великому ринку
3	Ціна продукту значно вища за ціни аналогів	Ціна продукту дещо вища за ціни аналогів	Ціна продукту приблизно дорівнює цінам аналогів	Ціна продукту дещо нижче за ціни аналогів	Ціна продукту значно нижче за ціни аналогів

## Продовження таблиці 5.1

4	Технічні та споживчі властивості продукту значно гірші, ніж в аналогів	Технічні та споживчі властивості продукту трохи гірші, ніж в аналогів	Технічні та споживчі властивості продукту на рівні аналогів	Технічні та споживчі властивості продукту трохи кращі, ніж в аналогів	Технічні та споживчі властивості продукту значно кращі, ніж в аналогів
5	Експлуатаційні витрати значно вищі, ніж в аналогів	Експлуатаційні витрати дещо вищі, ніж в аналогів	Експлуатаційні витрати на рівні експлуатаційних витрат аналогів	Експлуатаційні витрати трохи нижчі, ніж в аналогів	Експлуатаційні витрати значно нижчі, ніж в аналогів
Ринкові перспективи					
6	Ринок малий і не має позитивної динаміки	Ринок малий, але має позитивну динаміку	Середній ринок з позитивною динамікою	Великий стабільний ринок	Великий ринок з позитивною динамікою
7	Активна конкуренція великих компаній на	Активна конкуренція	Помірна конкуренція	Незначна конкуренція	Конкуренція немає
Практична здійсненність					
8	Відсутні фахівці як з технічної, так і з комерційної реалізації ідеї	Необхідно наймати фахівців або витратити значні кошти та час на навчання наявних фахівців	Необхідне незначне навчання фахівців та збільшення їх штату	Необхідне незначне навчання фахівців	Є фахівці з питань як з технічної, так і з комерційної реалізації ідеї
9	Потрібні значні фінансові ресурси, які відсутні. Джерела фінансування ідеї відсутні	Потрібні незначні фінансові ресурси. Джерела фінансування відсутні	Потрібні значні фінансові ресурси. Джерела фінансування є	Потрібні незначні фінансові ресурси. Джерела фінансування є	Не потребує додаткового фінансування
10	Необхідна розробка нових матеріалів	Потрібні матеріали, що використовуються у військово-промисловому комплексі	Потрібні дорогі матеріали	Потрібні досяжні та дешеві матеріали	Всі матеріали для реалізації ідеї відомі та давно використовуються у виробництві

## Продовження таблиці 5.1

11	Термін реалізації ідеї більший за 10 років	Термін реалізації ідеї більший за 5 років. Термін окупності інвестицій більше 10-ти років	Термін реалізації ідеї від 3-х до 5-ти років. Термін окупності інвестицій більше 5-ти років	Термін реалізації ідеї менше 3-х років. Термін окупності інвестицій від 3-х до 5-ти років	Термін реалізації ідеї менше 3-х років. Термін окупності інвестицій менше 3-х років
12	Необхідна розробка регламентних документів та отримання великої кількості дозвільних документів на виробництво та реалізацію продукту	Необхідно отримання великої кількості дозвільних документів на виробництво та реалізацію продукту, що вимагає значних коштів та часу	Процедура отримання дозвільних документів для виробництва та реалізації продукту вимагає незначних коштів та часу	Необхідно тільки повідомлення відповідним органам про виробництво та реалізацію продукту	Відсутні будь-які регламентні обмеження на виробництво та реалізацію продукту

Результати оцінювання науково-технічного рівня та комерційного потенціалу науково-технічної розробки потрібно звести до таблиці.

Таблиця 5.2 – Результати оцінювання науково-технічного рівня і комерційного потенціалу розробки експертами

Критерії	Експерт (ПІБ, посада)		
	1	2	3
	Бали:		
1. Технічна здійсненність концепції	5	5	5
2. Ринкові переваги (наявність аналогів)	4	5	4
3. Ринкові переваги (ціна продукту)	4	4	4
4. Ринкові переваги (технічні властивості)	5	4	5
5. Ринкові переваги (експлуатаційні витрати)	1	1	1
6. Ринкові перспективи (розмір ринку)	3	3	3
7. Ринкові перспективи (конкуренція)	1	2	1
8. Практична здійсненність (наявність фахівців)	3	3	3
9. Практична здійсненність (наявність фінансів)	3	4	4
10. Практична здійсненність (необхідність нових матеріалів)	4	4	4
11. Практична здійсненність (термін реалізації)	4	4	4
12. Практична здійсненність (розробка документів)	4	4	5
Сума балів	41	43	43
Середньоарифметична сума балів $СБ_c$	42,3		

За результатами розрахунків, наведених в таблиці 5.2, зробимо висновок щодо науково-технічного рівня і рівня комерційного потенціалу розробки. При цьому використаємо рекомендації, наведені в табл. 5.3 [28].

Таблиця 5.3 – Науково-технічні рівні та комерційні потенціали розробки

Середньоарифметична сума балів СБ розрахована на основі висновків експертів	Науково-технічний рівень та комерційний потенціал розробки
41...48	Високий
31...40	Вище середнього
21...30	Середній
11...20	Нижче середнього
0...10	Низький

Згідно проведених досліджень рівень комерційного потенціалу розробки за темою «Інформаційна технологія прогнозування успішності кінофільму» становить 42,3 бала, що, відповідно до таблиці 4.3, свідчить про комерційну важливість проведення даних досліджень (рівень комерційного потенціалу розробки високий).

Перед переглядом кінофільмів компаніями має забезпечуватись інформування глядачів про кінофільми, їхній бюджет, прогнозований успіх та касовий збір. Перед тим, як замовляти білети до кінотеатру користувач може перевірити, чи потрібно йому витратити гроші та час на прем'єру конкретного фільму, чи ні. Прогнозування успіху – це саме те, чого потребує користувач, який буде використовувати розроблений веб-сервіс. Основні відмінності від конкурентів полягають в наступному: 1) конкурентів у відкритому доступі майже немає; 2) покращений інтерфейс для користувача; 3) збільшена точність прогнозування; 4) забезпечена можливість додавання власних даних про кінофільм, та отримувати на основі цього прогнозування, в той час як у конкурентів можна переглянути тільки вже існуючі кінофільми; 5) програмний продукт також відображає, який параметр найбільше повипливав

на прогноз.

Використані програмні засоби:

IDE - PyCharm, мова програмування Python, Google Cloud Platform - для розгортання серверної частини та бази даних, Firebase - для хостингу клієнтської частини.

## 5.2 Розрахунок узагальненого коефіцієнта якості розробки

Окрім комерційного аудиту розробки доцільно також розглянути технічний рівень якості розробки, розглянувши її основні технічні показники. Ці показники по-різному впливають на загальну якість проектної розробки.

Узагальнений коефіцієнт якості ( $B_n$ ) для нового технічного рішення розрахуємо за формулою [29]:

$$B_n = \sum_{i=1}^k \alpha_i \cdot \beta_i, \quad (5.1)$$

де  $k$  – кількість найбільш важливих технічних показників, які впливають на якість нового технічного рішення;

$\alpha_i$  – коефіцієнт, який враховує питому вагу  $i$ -го технічного показника в загальній якості розробки. Коефіцієнт  $\alpha_i$  визначається експертним

шляхом і при цьому має виконуватись умова  $\sum_{i=1}^k \alpha_i = 1$ ;

$\beta_i$  – відносне значення  $i$ -го технічного показника якості нової розробки.

Відносні значення  $\beta_i$  для різних випадків розраховуємо за такими формулами:

- для показників, зростання яких вказує на підвищення в лінійній залежності якості нової розробки:

$$\beta_i = \frac{I_{ni}}{I_{na}}, \quad (5.2)$$

де  $I_{ni}$  та  $I_{na}$  – чисельні значення конкретного  $i$ -го технічного показника



якості відповідно для нової розробки та аналога;

- для показників, зростання яких вказує на погіршення в лінійній залежності якості нової розробки:

$$\beta_i = \frac{I_{ai}}{I_{ni}}; \quad (5.3)$$

Використовуючи наведені залежності можемо проаналізувати та порівняти техніко-економічні характеристики аналогу та розробки на основі отриманих наявних та проектних показників, а результати порівняння зведемо до таблиці 5.4.

Таблиця 5.4 – Порівняння основних параметрів розробки та аналога

Показники (параметри)	Одиниця вимірювання	Аналог	Проектований пристрій	Відношення параметрів нової розробки до аналога	Питома вага показника
Задіяні операційні системи	-	2	2	1	0,15
Дружність інтерфейсу	бал	6	8	1,33	0,1
Кількість базових характеристик розпізнавання	шт.	32	64	2	0,3
Точність обробки прогнозу	%	85	94	1,11	0,25
Функціональність	шт.	3	6	2	0,2

Узагальнений коефіцієнт якості ( $B_n$ ) для нового технічного рішення складе:

$$B_n = \sum_{i=1}^k \alpha_i \cdot \beta_i = 1 \cdot 0,15 + 1,33 \cdot 0,1 + 2 \cdot 0,3 + 1,11 \cdot 0,25 + 2 \cdot 0,2 = 1,56.$$

Отже за технічними параметрами, згідно узагальненого коефіцієнту якості розробки, науково-технічна розробка переважає існуючі аналоги приблизно в 1,56 рази.

### 5.3 Розрахунок витрат на проведення науково-дослідної роботи

Витрати, пов'язані з проведенням науково-дослідної роботи на тему «Інформаційна технологія прогнозування успішності кінофільму», під час планування, обліку і калькулювання собівартості науково-дослідної роботи групуємо за відповідними статтями.

#### 5.3.1 Витрати на оплату праці

До статті «Витрати на оплату праці» належать витрати на виплату основної та додаткової заробітної плати керівникам відділів, лабораторій, секторів і груп, науковим, інженерно-технічним працівникам, конструкторам, технологам, креслярам, копіювальникам, лаборантам, робітникам, студентам, аспірантам та іншим працівникам, безпосередньо зайнятим виконанням конкретної теми, обчисленої за посадовими окладами, відрядними розцінками, тарифними ставками згідно з чинними в організаціях системами оплати праці.

#### Основна заробітна плата дослідників

Витрати на основну заробітну плату дослідників ( $Z_o$ ) розраховуємо у відповідності до посадових окладів працівників, за формулою [28]:

$$Z_o = \sum_{i=1}^k \frac{M_{ni} \cdot t_i}{T_p}, \quad (5.4)$$

де  $k$  – кількість посад дослідників залучених до процесу досліджень;

$M_{ni}$  – місячний посадовий оклад конкретного дослідника, грн;

$t_i$  – число днів роботи конкретного дослідника, дн.;

$T_p$  – середнє число робочих днів в місяці,  $T_p=21$  дні.

$$Z_o = 14830,00 \cdot 34 / 21 = 24010,48 \text{ грн.}$$

Проведені розрахунки зведемо до таблиці 5.5.

Таблиця 5.5 – Витрати на заробітну плату дослідників

Найменування посади	Місячний посадовий оклад, грн	Оплата за робочий день, грн	Число днів роботи	Витрати на заробітну плату, грн
Керівник проекту	14830,00	706,19	34	24010,48
Консультант (менеджер сфери прокату фільмів)	12720,00	605,71	15	9085,71
Інженер-програміст	12100,00	576,19	20	11523,81
Фахівець з аналітично-математичних досліджень	12000,00	571,43	20	11428,57
Консультант-аналітик цифрових обчислюваних систем	13800,00	657,14	15	9857,14
Лаборант	6800,00	323,81	20	6476,19
Всього				72381,90

### Основна заробітна плата робітників

Витрати на основну заробітну плату робітників ( $Z_p$ ) за відповідними найменуваннями робіт НДР на тему «Інформаційна технологія прогнозування успішності кінофільму» розраховуємо за формулою:

$$Z_p = \sum_{i=1}^n C_i \cdot t_i, \quad (5.5)$$

де  $C_i$  – погодинна тарифна ставка робітника відповідного розряду, за виконану відповідну роботу, грн/год;

$t_i$  – час роботи робітника при виконанні визначеної роботи, год.

Погодинну тарифну ставку робітника відповідного розряду  $C_i$  можна визначити за формулою:

$$C_i = \frac{M_M \cdot K_i \cdot K_c}{T_p \cdot t_{зм}}, \quad (5.6)$$

де  $M_M$  – розмір прожиткового мінімуму працездатної особи, або мінімальної місячної заробітної плати (в залежності від діючого законодавства), прийmemo  $M_M=6700,00$  грн;

$K_i$  – коефіцієнт міжкваліфікаційного співвідношення для встановлення тарифної ставки робітнику відповідного розряду (табл. Б.2, додаток Б) [28];

$K_c$  – мінімальний коефіцієнт співвідношень місячних тарифних ставок робітників першого розряду з нормальними умовами праці виробничих об'єднань і підприємств до законодавчо встановленого розміру мінімальної заробітної плати.

$T_p$  – середнє число робочих днів в місяці, приблизно  $T_p = 21$  дн;

$t_{зм}$  – тривалість зміни, год.

$$C_l = 6700,00 \cdot 1,10 \cdot 1,7 / (21 \cdot 8) = 74,58 \text{ грн.}$$

$$З_{pl} = 74,58 \cdot 8,00 = 596,62 \text{ грн.}$$

Таблиця 5.6 – Величина витрат на основну заробітну плату робітників

Найменування робіт	Тривалість роботи, год	Розряд роботи	Тарифний коефіцієнт	Погодинна тарифна ставка, грн	Величина оплати на робітника грн
Встановлення допоміжного офісного обладнання	8,00	2	1,10	74,58	596,62
Монтаж робочого місця розробника системи прогнозування	12,00	2	1,10	74,58	894,93
Інсталяція програмного забезпечення	5,00	5	1,70	115,26	576,28
Встановлення цифрових обчислювальних систем	3,00	4	1,50	101,70	305,09
Відлагодження інтерполяційних модулів	7,00	5	1,70	115,26	806,79
Тренування цифрової експериментальної моделі	4,50	4	1,50	101,70	457,63
Формування бази даних прогнозного аналізу	16,00	3	1,35	91,53	1464,43
Інші допоміжні роботи	10,00	3	1,35	91,53	915,27
Всього					6017,04

Додаткова заробітна плата дослідників та робітників

Додаткову заробітну плату розраховуємо як 10 ... 12% від суми основної заробітної плати дослідників та робітників за формулою:

$$З_{доd} = (З_o + З_p) \cdot \frac{H_{доd}}{100\%}, \quad (5.7)$$

де  $H_{доd}$  – норма нарахування додаткової заробітної плати. Прийmemo 10%.

$$Z_{\text{од}} = (72381,90 + 6017,04) \cdot 10 / 100\% = 7839,89 \text{ грн.}$$

### 5.3.2 Відрахування на соціальні заходи

Нарахування на заробітну плату дослідників та робітників розраховуємо як 22% від суми основної та додаткової заробітної плати дослідників і робітників за формулою:

$$Z_n = (Z_o + Z_p + Z_{\text{од}}) \cdot \frac{H_{\text{зн}}}{100\%} \quad (5.8)$$

де  $H_{\text{зн}}$  – норма нарахування на заробітну плату. Приймаємо 22%.

$$Z_n = (72381,90 + 6017,04 + 7839,89) \cdot 22 / 100\% = 18972,54 \text{ грн.}$$

### 5.3.3 Сировина та матеріали

До статті «Сировина та матеріали» належать витрати на сировину, основні та допоміжні матеріали, інструменти, пристрої та інші засоби і предмети праці, які придбані у сторонніх підприємств, установ і організацій та витрачені на проведення досліджень за темою «Інформаційна технологія прогнозування успішності кінофільму».

Витрати на матеріали ( $M$ ), у вартісному вираженні розраховуються окремо по кожному виду матеріалів за формулою:

$$M = \sum_{j=1}^n H_j \cdot C_j \cdot K_j - \sum_{j=1}^n B_j \cdot C_{\text{в}j}, \quad (5.9)$$

де  $H_j$  – норма витрат матеріалу  $j$ -го найменування, кг;

$n$  – кількість видів матеріалів;

$C_j$  – вартість матеріалу  $j$ -го найменування, грн/кг;

$K_j$  – коефіцієнт транспортних витрат, ( $K_j = 1,1 \dots 1,15$ );

$B_j$  – маса відходів  $j$ -го найменування, кг;

$C_{\text{в}j}$  – вартість відходів  $j$ -го найменування, грн/кг.

$$M_1 = 3,00 \cdot 265,00 \cdot 1,1 - 0,000 \cdot 0,00 = 874,50 \text{ грн.}$$

Проведені розрахунки зведемо до таблиці 5.7.

Таблиця 5.7 – Витрати на матеріали

Найменування матеріалу, марка, тип, сорт	Ціна за 1 кг, грн	Норма витрат, кг	Величина відходів, кг	Ціна відходів, грн/кг	Вартість витраченого матеріалу, грн
Папір канцелярський офісний ECONOMIC (A4-500)	265,00	3,00	-	-	874,50
Папір для заміток ECONOMIC (A5)-60	162,00	4,00	-	-	712,80
Начиння канцелярське DATUM FX	22,00	3,00	-	-	72,60
Органайзер офісний DATUM Office	150,00	3,00	-	-	495,00
Картридж для принтера HP-210A	1250,00	2,00	-	-	2750,00
Диск оптичний VEKO-10 (CD-R)	22,00	1,00	-	-	24,20
Диск оптичний VEKO-W (CD-RW)	23,00	1,00	-	-	25,30
FLASH-пам'ять Kingstar (32 ГБ) Class 10	340,00	1,00	-	-	372,00
FLASH-пам'ять Kingstar (64 ГБ) Class 10 A	680,00	1,000	-	-	748,00
Всього					6072,10

#### 5.3.4 Розрахунок витрат на комплектуючі

Витрати на комплектуючі ( $K_6$ ), які використовують при проведенні НДР на тему «Інформаційна технологія прогнозування успішності кінофільму» відсутні.

#### 5.3.5 Спецустаткування для наукових (експериментальних) робіт

До статті «Спецустаткування для наукових (експериментальних) робіт» належать витрати на виготовлення та придбання спецустаткування необхідного для проведення досліджень, також витрати на їх проектування, виготовлення, транспортування, монтаж та встановлення.

Балансову вартість спекустаткування розраховуємо за формулою:

$$B_{\text{спец}} = \sum_{i=1}^k C_i \cdot C_{\text{пр.}i} \cdot K_i, \quad (5.10)$$

де  $C_i$  – ціна придбання одиниці спекустаткування даного виду, марки, грн;

$C_{\text{пр.}i}$  – кількість одиниць устаткування відповідного найменування, які придбані для проведення досліджень, шт.;

$K_i$  – коефіцієнт, що враховує доставку, монтаж, налагодження устаткування тощо, ( $K_i = 1,10 \dots 1,12$ );

$k$  – кількість найменувань устаткування.

$$B_{\text{спец}} = 40250,00 \cdot 1 \cdot 1,11 = 44677,50 \text{ грн.}$$

Отримані результати зведемо до таблиці 5.8.

Таблиця 5.8 – Витрати на придбання спекустаткування по кожному виду

Найменування устаткування	Кількість, шт	Ціна за одиницю, грн	Вартість, грн
Мультимедійна проекційна система	1	40250,00	44677,50
Всього			44677,50

### 5.3.6 Програмне забезпечення для наукових (експериментальних) робіт

До статті «Програмне забезпечення для наукових (експериментальних) робіт» належать витрати на розробку та придбання спеціальних програмних засобів і програмного забезпечення, (програм, алгоритмів, баз даних) необхідних для проведення досліджень, також витрати на їх проектування, формування та встановлення.

Балансову вартість програмного забезпечення розраховуємо за формулою:

$$B_{\text{прог}} = \sum_{i=1}^k C_{\text{инрг}} \cdot C_{\text{прог.}i} \cdot K_i, \quad (5.11)$$

де  $C_{\text{инрг}}$  – ціна придбання одиниці програмного засобу даного виду, грн;

$C_{прз.i}$  – кількість одиниць програмного забезпечення відповідного найменування, які придбані для проведення досліджень, шт.;

$K_i$  – коефіцієнт, що враховує інсталяцію, налагодження програмного засобу тощо, ( $K_i = 1,10...1,12$ );

$k$  – кількість найменувань програмних засобів.

$$B_{прз} = 8410,00 \cdot 1 \cdot 1,1 = 9251,00 \text{ грн.}$$

Отримані результати зведемо до таблиці 5.9.

Таблиця 5.9– Витрати на придбання програмних засобів по кожному виду

Найменування програмного засобу	Кількість, шт	Ціна за одиницю, грн	Вартість, грн
ОС Windows	1	8410,00	9251,00
Прикладний пакет Microsoft Office	1	7790,00	8569,00
Програмний засіб IDE – PyCharm	1	6580,00	7238,00
Google Cloud Platform	1	6830,00	7513,00
Firebase	1	5935,00	6528,50
Всього			39099,50

### 5.3.7 Амортизація обладнання, програмних засобів та приміщень

В спрощеному вигляді амортизаційні відрахування по кожному виду обладнання, приміщень та програмному забезпеченню тощо, розраховуємо з використанням прямолінійного методу амортизації за формулою:

$$A_{обл} = \frac{Ц_б}{T_г} \cdot \frac{t_{вик}}{12}, \quad (5.12)$$

де  $Ц_б$  – балансова вартість обладнання, програмних засобів, приміщень тощо, які використовувались для проведення досліджень, грн;

$t_{вик}$  – термін використання обладнання, програмних засобів, приміщень під час досліджень, місяців;

$T_г$  – строк корисного використання обладнання, програмних засобів, приміщень тощо, років.

$$A_{обл} = (28300,00 \cdot 2) / (3 \cdot 12) = 1572,22 \text{ грн.}$$

Проведені розрахунки зведемо до таблиці 5.10.



Таблиця 5.10 – Амортизаційні відрахування по кожному виду обладнання

Найменування обладнання	Балансова вартість, грн	Строк корисного використання, років	Термін використання обладнання, місяців	Амортизаційні відрахування, грн
Програмно-аналітичний комплекс	28300,00	3	2	1572,22
Графічно-обчислювальний комплекс обробки даних	26400,00	4	2	1100,00
Програмні засоби реалізації Python	6890,00	2	2	574,17
Обладнання виводу інформації	10250,00	5	2	341,67
Місце оператора спеціалізоване	9100,00	4	2	379,17
Офісна оргтехніка	9600,00	5	2	320,00
Дослідницька лабораторія	300000,00	25	2	2000,00
Всього				6287,22

### 5.3.8 Паливо та енергія для науково-виробничих цілей

Витрати на силову електроенергію ( $B_e$ ) розраховуємо за формулою:

$$B_e = \sum_{i=1}^n \frac{W_{yi} \cdot t_i \cdot C_e \cdot K_{eni}}{\eta_i}, \quad (5.13)$$

де  $W_{yi}$  – встановлена потужність обладнання на визначеному етапі розробки, кВт;

$t_i$  – тривалість роботи обладнання на етапі дослідження, год;

$C_e$  – вартість 1 кВт-години електроенергії, грн; (вартість електроенергії визначається за даними енергопостачальної компанії), прийmemo  $C_e = 6,20$  грн;

$K_{eni}$  – коефіцієнт, що враховує використання потужності,  $K_{eni} < 1$ ;

$\eta_i$  – коефіцієнт корисної дії обладнання,  $\eta_i < 1$ .

$$B_e = 0,28 \cdot 240,0 \cdot 6,20 \cdot 0,95 / 0,97 = 416,64 \text{ грн.}$$

Проведені розрахунки зведемо до таблиці 5.11.

Таблиця 5.11 – Витрати на електроенергію

Найменування обладнання	Встановлена потужність, кВт	Тривалість роботи, год	Сума, грн
Програмно-аналітичний комплекс	0,28	240,0	416,64
Графічно-обчислювальний комплекс обробки даних	0,32	200,0	396,80
Мультимедійна проекційна система	0,32	200,0	396,80
Обладнання виводу інформації	0,35	40,0	86,80
Місце оператора спеціалізоване	0,46	200,0	570,40
Офісна оргтехніка	0,80	24,0	119,04
Всього			1986,48

### 5.3.9 Службові відрядження

До статті «Службові відрядження» дослідної роботи на тему «Інформаційна технологія прогнозування успішності кінофільму» належать витрати на відрядження штатних працівників, працівників організацій, які працюють за договорами цивільно-правового характеру, аспірантів, зайнятих розробленням досліджень, відрядження, пов'язані з проведенням випробувань машин та приладів, а також витрати на відрядження на наукові з'їзди, конференції, наради, пов'язані з виконанням конкретних досліджень.

Витрати за статтею «Службові відрядження» розраховуємо як 20...25% від суми основної заробітної плати дослідників та робітників за формулою:

$$B_{cb} = (Z_o + Z_p) \cdot \frac{H_{cb}}{100\%}, \quad (5.15)$$

де  $H_{cb}$  – норма нарахування за статтею «Службові відрядження», прийmemo  $H_{cb} = 25\%$ .

$$B_{cb} = (72381,90 + 6017,04) \cdot 25 / 100\% = 19599,74 \text{ грн.}$$

5.3.10 Витрати на роботи, які виконують сторонні підприємства, установи і організації

Витрати за статтею «Витрати на роботи, які виконують сторонні підприємства, установи і організації» розраховуємо як 30...45% від суми основної заробітної плати дослідників та робітників за формулою:

$$B_{cn} = (Z_o + Z_p) \cdot \frac{H_{cn}}{100\%}, \quad (5.16)$$

де  $H_{cn}$  – норма нарахування за статтею «Витрати на роботи, які виконують сторонні підприємства, установи і організації», прийmemo  $H_{cn} = 40\%$ .

$$B_{cn} = (72381,90 + 6017,04) \cdot 40 / 100\% = 31359,58 \text{ грн.}$$

### 5.3.11 Інші витрати

До статті «Інші витрати» належать витрати, які не знайшли відображення у зазначених статтях витрат і можуть бути віднесені безпосередньо на собівартість досліджень за прямими ознаками.

Витрати за статтею «Інші витрати» розраховуємо як 50...100% від суми основної заробітної плати дослідників та робітників за формулою:

$$I_g = (Z_o + Z_p) \cdot \frac{H_{ig}}{100\%}, \quad (5.17)$$

де  $H_{ig}$  – норма нарахування за статтею «Інші витрати», прийmemo  $H_{ig} = 100\%$ .

$$I_g = (72381,90 + 6017,04) \cdot 100 / 100\% = 78398,94 \text{ грн.}$$

### 5.3.12 Накладні (загальновиробничі) витрати

До статті «Накладні (загальновиробничі) витрати» належать: витрати, пов'язані з управлінням організацією; витрати на винахідництво та раціоналізацію; витрати на підготовку (перепідготовку) та навчання кадрів; витрати, пов'язані з набором робочої сили; витрати на оплату послуг банків;

витрати, пов'язані з освоєнням виробництва продукції; витрати на науково-технічну інформацію та рекламу та ін.

Витрати за статтею «Накладні (загальновиробничі) витрати» розраховуємо як 100...150% від суми основної заробітної плати дослідників та робітників за формулою:

$$B_{нзв} = (Z_o + Z_p) \cdot \frac{H_{нзв}}{100\%}, \quad (5.18)$$

де  $H_{нзв}$  – норма нарахування за статтею «Накладні (загальновиробничі) витрати», прийmemo  $H_{нзв} = 150\%$ .

$$B_{нзв} = (72381,90 + 6017,04) \cdot 150 / 100\% = 117598,42 \text{ грн.}$$

Витрати на проведення науково-дослідної роботи на тему «Інформаційна технологія прогнозування успішності кінофільму» розраховуємо як суму всіх попередніх статей витрат за формулою:

$$B_{заг} = Z_o + Z_p + Z_{доо} + Z_n + M + K_v + B_{спец} + B_{прз} + A_{обл} + B_e + B_{св} + B_{сп} + I_v + B_{нзв}. \quad (5.19)$$

$$B_{заг} = 72381,90 + 6017,04 + 7839,89 + 18972,54432 + 6072,10 + 0,00 + 44677,50 + 39099,50 + 6287,22 + 1986,48 + 19599,74 + 31359,58 + 78398,94 + 117598,42 = 450290,86 \text{ грн.}$$

Загальні витрати  $ЗВ$  на завершення науково-дослідної (науково-технічної) роботи та оформлення її результатів розраховується за формулою:

$$ЗВ = \frac{B_{заг}}{\eta}, \quad (5.20)$$

де  $\eta$  - коефіцієнт, який характеризує етап (стадію) виконання науково-дослідної роботи, прийmemo  $\eta = 0,85$ .

$$ЗВ = 450290,86 / 0,85 = 529753,95 \text{ грн.}$$

#### 5.4 Розрахунок економічної ефективності науково-технічної розробки при її можливій комерціалізації потенційним інвестором

В ринкових умовах узагальнюючим позитивним результатом, що його може отримати потенційний інвестор від можливого впровадження результатів тієї чи іншої науково-технічної розробки, є збільшення у потенційного інвестора величини чистого прибутку.

Результати дослідження проведені за темою «Інформаційна технологія прогнозування успішності кінофільму» передбачають комерціалізацію протягом 4-х років реалізації на ринку.

В цьому випадку основу майбутнього економічного ефекту будуть формувати:

$\Delta N$  – збільшення кількості споживачів яким надається відповідна інформаційна послуга у періоди часу, що аналізуються;

Показник	1-й рік	2-й рік	3-й рік	4-й рік
Збільшення кількості споживачів, осіб	1000	1500	1100	650

$N$  – кількість споживачів яким надавалась відповідна інформаційна послуга у році до впровадження результатів нової науково-технічної розробки, прийmemo 5800 осіб;

$C_o$  – вартість послуги у році до впровадження інформаційної системи, прийmemo 3600,00 грн;

$\pm \Delta C_o$  – зміна вартості послуги від впровадження результатів, прийmemo 1897,80 грн.

Можливе збільшення чистого прибутку у потенційного інвестора  $\Delta \Pi_i$  для кожного із 4-х років, протягом яких очікується отримання позитивних результатів від можливого впровадження та комерціалізації науково-технічної розробки, розраховуємо за формулою [28]:

$$\Delta\Pi_i = (\pm\Delta C_o \cdot N + C_o \cdot \Delta N)_i \cdot \lambda \cdot \rho \cdot \left(1 - \frac{\vartheta}{100}\right), \quad (5.21)$$

де  $\lambda$  – коефіцієнт, який враховує сплату потенційним інвестором податку на додану вартість. У 2022 році ставка податку на додану вартість складає 20%, а коефіцієнт  $\lambda = 0,8333$ ;

$\rho$  – коефіцієнт, який враховує рентабельність інноваційного продукту).

Прийmemo  $\rho = 30\%$ ;

$\vartheta$  – ставка податку на прибуток, який має сплачувати потенційний інвестор, у 2022 році  $\vartheta = 18\%$ ;

Збільшення чистого прибутку 1-го року:

$$\Delta\Pi_1 = (1897,80 \cdot 5800,00 + 5497,80 \cdot 1000) \cdot 0,83 \cdot 0,3 \cdot (1 - 0,18/100\%) = 3369999,07 \text{ грн.}$$

Збільшення чистого прибутку 2-го року:

$$\Delta\Pi_2 = (1897,80 \cdot 5800,00 + 5497,80 \cdot 2500) \cdot 0,83 \cdot 0,3 \cdot (1 - 0,18/100\%) = 5053810,27 \text{ грн.}$$

Збільшення чистого прибутку 3-го року:

$$\Delta\Pi_3 = (1897,80 \cdot 5800,00 + 5497,80 \cdot 3600) \cdot 0,83 \cdot 0,3 \cdot (1 - 0,18/100\%) = 6288605,16 \text{ грн.}$$

Збільшення чистого прибутку 4-го року:

$$\Delta\Pi_4 = (1897,80 \cdot 5800,00 + 5497,80 \cdot 4250) \cdot 0,83 \cdot 0,3 \cdot (1 - 0,18/100\%) = 7018256,68 \text{ грн.}$$

Приведена вартість збільшення всіх чистих прибутків  $ПП$ , що їх може отримати потенційний інвестор від можливого впровадження та комерціалізації науково-технічної розробки:

$$ПП = \sum_{i=1}^T \frac{\Delta\Pi_i}{(1 + \tau)^t}, \quad (5.22)$$

де  $\Delta\Pi_i$  – збільшення чистого прибутку у кожному з років, протягом яких виявляються результати впровадження науково-технічної розробки, грн;

$T$  – період часу, протягом якого очікується отримання позитивних результатів від впровадження та комерціалізації науково-технічної розробки, роки;

$\tau$  – ставка дисконтування, за яку можна взяти щорічний прогнозований рівень інфляції в країні,  $\tau = 0,18$ ;

$t$  – період часу (в роках) від моменту початку впровадження науково-технічної розробки до моменту отримання потенційним інвестором додаткових чистих прибутків у цьому році.

$$\begin{aligned} III &= 3369999,07/(1+0,18)^1 + 5053810,27/(1+0,18)^2 + 6288605,16/(1+0,18)^3 + \\ &+ 7018256,68/(1+0,18)^4 = 2855931,41 + 3629567,85 + 3827439,25 + 3619938,72 = \\ &= 13932877,23 \text{ грн.} \end{aligned}$$

**Величина початкових інвестицій  $PV$ , які потенційний інвестор має вкласти для впровадження і комерціалізації науково-технічної розробки:**

$$PV = k_{инв} \cdot 3B, \quad (5.23)$$

де  $k_{инв}$  – коефіцієнт, що враховує витрати інвестора на впровадження науково-технічної розробки та її комерціалізацію, приймаємо  $k_{инв} = 1,5$ ;

$3B$  – загальні витрати на проведення науково-технічної розробки та оформлення її результатів, приймаємо 529753,95 грн.

$$PV = k_{инв} \cdot 3B = 1,5 \cdot 529753,95 = 794630,92 \text{ грн.}$$

Абсолютний економічний ефект  $E_{абс}$  для потенційного інвестора від можливого впровадження та комерціалізації науково-технічної розробки становитиме:

$$E_{абс} = III - PV \quad (5.24)$$

де  $ПП$  – приведена вартість зростання всіх чистих прибутків від можливого впровадження та комерціалізації науково-технічної розробки, 13932877,23 грн;

$PV$  – теперішня вартість початкових інвестицій, 794630,92 грн.

$E_{абс} = ПП - PV = 13932877,23 - 794630,92 = 13138246,30$  грн.

Внутрішня економічна прибутковість інвестицій  $E_ε$ , які можуть бути вкладені потенційним інвестором у впровадження та комерціалізацію науково-технічної розробки:

$$E_ε = T_{жс} \sqrt[4]{1 + \frac{E_{абс}}{PV}} - 1, \quad (5.25)$$

де  $E_{абс}$  – абсолютний економічний ефект вкладених інвестицій, 13138246,30 грн;

$PV$  – теперішня вартість початкових інвестицій, 794630,92 грн;

$T_{жс}$  – життєвий цикл науково-технічної розробки, тобто час від початку її розробки до закінчення отримання позитивних результатів від її впровадження, 4 роки.

$$E_ε = T_{жс} \sqrt[4]{1 + \frac{E_{абс}}{PV}} - 1 = (1 + 13138246,30/794630,92)^{1/4} = 1,05.$$

Мінімальна внутрішня економічна прибутковість вкладених інвестицій

$\tau_{мін}$

$$\tau_{мін} = d + f, \quad (5.26)$$

де  $d$  – середньозважена ставка за депозитними операціями в комерційних банках; в 2022 році в Україні  $d = 0,12$ ;

$f$  – показник, що характеризує ризикованість вкладення інвестицій, приймемо 0,3.



$\tau_{\min} = 0,12 + 0,3 = 0,42 < 1,05$  свідчить про те, що внутрішня економічна прибутковість інвестицій  $E_g$ , які можуть бути вкладені потенційним інвестором у впровадження та комерціалізацію науково-технічної розробки вища мінімальної внутрішньої прибутковості. Тобто інвестувати в науково-дослідну роботу за темою «Інформаційна технологія прогнозування успішності кінофільму» доцільно.

Період окупності інвестицій  $T_{ок}$  які можуть бути вкладені потенційним інвестором у впровадження та комерціалізацію науково-технічної розробки:

$$T_{ок} = \frac{1}{E_g}, \quad (5.27)$$

де  $E_g$  – внутрішня економічна прибутковість вкладених інвестицій;  $T_{ок} = 1 / 1,05 = 0,96$  р;  $T_{ок} < 3$ -х років, що свідчить про комерційну привабливість науково-технічної розробки і може спонукати потенційного інвестора профінансувати впровадження даної розробки та виведення її на ринок.

#### 5.4 Висновок до розділу 5

Згідно проведених досліджень рівень комерційного потенціалу розробки за темою «Інформаційна технологія прогнозування успішності кінофільму» становить 42,3 бала, що, свідчить про комерційну важливість проведення даних досліджень (рівень комерційного потенціалу розробки високий).

При оцінюванні за технічними параметрами, згідно узагальненого коефіцієнту якості розробки, науково-технічна розробка переважає існуючі аналоги приблизно в 1,56 рази.

Також термін окупності становить 0,96 р., що менше 3-х років, що свідчить про комерційну привабливість науково-технічної розробки і може спонукати потенційного інвестора профінансувати впровадження даної розробки та виведення її на ринок.

Отже можна зробити висновок про доцільність проведення науково-дослідної роботи за темою «Інформаційна технологія прогнозування успішності кінофільму».

## ВИСНОВКИ

При виконанні магістерської кваліфікаційної роботи розв'язано задачу розробки інформаційної технології та програмного забезпечення розв'язання задачі прогнозування успішності кінофільму основі власної структури регресійної нейронної мережі.

В першому розділі проведено аналіз предметної області прогнозування успішності кінофільму, а саме – сформульовано постановку задачі, проведено огляд відомих методів розв'язання задачі прогнозування успішності кінофільму. Виходячи з аналізів предметної області, дало можливість визначити основні вимоги до системи прогнозування успішності кінофільму. Найперспективнішим було визначено метод на основі штучних нейронних мереж. Проведено аналіз існуючих методів, так виявлено, що можна підвищити точність прогнозування успішності кінофільмів.

У другому розділі магістерської кваліфікаційної роботи було проаналізовано ознаки кінофільмів, та визначено, ті які найбільше впливають на прибуток кінофільму. Запропоновано власну структура нейронної мережі на базі регресійної моделі. Кінцева модель складає: шар вхідних даних, два прихований шар, та один вихідний шар з 1 нейроном, для отримання скалярної величини на виході.

У третьому розділі наведено основні кроки роботи інформаційної технології прогнозування успішності кінофільму. Обрано клієнт-серверну архітектуру для програмного забезпечення прогнозування успішності кінофільму. Розбито програмний засіб на модулі, а саме: модуль взаємодії з базою даних, модуль імпортування кінофільмів, модуль редагування записів в базі даних, модуль прогнозування успішності кінофільму. Описано роботу та алгоритм кожного модуля. Спроектовано базу даних для програмного забезпечення прогнозування успішності кінофільмів.

У четвертому розділі обґрунтовано вибір мови програмування *Python* для серверної частини, та *Dart* – для клієнтської частини. У результаті було

розроблено програмне забезпечення прогнозування успішності кінофільму.

Було проведено тестування програми нейромережевого проказування успішності кінофільму. Аналіз результатів тестування показує, що точність прогнозування складає 94%, у той час коли точність аналогів складає 85% при навчанні на однаковій навчальній множині у кількість 5000 кінофільмів. Тобто мета магістерської кваліфікаційної роботи досягнута – точність прогнозування успішності кінофільму підвищена на 9%.

У п'ятому розділі було виконано розрахунок витрат на розробку та виготовлення нового технічного рішення. При оцінюванні за технічними параметрами, згідно узагальненого коефіцієнту якості розробки, науково-технічна розробка переважає існуючі аналоги приблизно в 1,56 рази. Термін окупності становить 0,96 р., що менше 3-х років, що свідчить про комерційну привабливість науково-технічної розробки і може спонукати потенційного інвестора профінансувати впровадження даної розробки та виведення її на ринок.

## ПЕРЕЛІК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. В. М. Борисюк, С. В. Барабан, «Виділення ознак кінофільмів, що впливають на успішність кінофільму», в Матеріали конференції «Молодь в науці: дослідження, проблеми, перспективи (МН-2023)», Вінниця, 2023, [Електронний ресурс] - режим доступу: <https://conferences.vntu.edu.ua/index.php/mn/mn2023/paper/view/16839/1404>.
2. L. Zhang, J. Luo, and S. Yang, “Forecasting box office revenue of movies with bp neural network,” Expert Systems with Applications, 2009.
3. V. Biramane, H. Kulkarni, A. Bhave, and P. Kosamkar, “Relationships between classical factors, social factors and box office collections,” 2016.
4. R. Masrury, M. Saputra, A. Alamsyah, and M. Primantari, “A comparative study of hollywood movie successfulness prediction model,” 2019.
5. S. S. Mohanbir, J. Eliashberg, «A Parsimonious Model for Forecasting Gross Box-Office Revenues of Motion Pictures», 1996.
6. J. Eliashberg et al, “The motion picture industry: critical issues in practice, current research, and new research directions”, Marketing Science, 2006.
7. Gagniuc, A. Paul, «Markov Chains: From Theory to Implementation and Experimentation», John Wiley & Sons, 2017. - 1–235 с.
8. D. A. Edwards, R. Buckmire, «A differential equation model of North American cinematic box-office dynamics», IMA Journal of management Mathematics, 2001.
9. R. Parimi and D. Caragea, “Pre-release box-office success prediction for motion pictures,” 2013.
10. L. Bottou, O. Bousquet, «The Tradeoffs of Large-Scale Learning», Cambridge: MIT Press, 2012. - 51–368 с.
11. Офіційна сторінка TMDb [Електронний ресурс] – режим доступу: <https://www.themoviedb.org/>.

12. Word cloud [Електронний ресурс] – режим доступу: [https://en.wikipedia.org/wiki/Tag\\_cloud](https://en.wikipedia.org/wiki/Tag_cloud).
13. *Wordcloud* документація [Електронний ресурс] – режим доступу: [http://amueller.github.io/word\\_cloud/](http://amueller.github.io/word_cloud/).
14. A Rajaraman, J. D. Ulman, «Mining of Massive Datasets», 2011. – 1-17 с.
15. Офіційна сторінка *eli5* [Електронний ресурс] – режим доступу: <https://eli5.readthedocs.io/en/latest/>.
16. Коефіцієнт детермінації [Електронний ресурс] – режим доступу: [https://en.wikipedia.org/wiki/Coefficient\\_of\\_determination](https://en.wikipedia.org/wiki/Coefficient_of_determination)
17. Середня квадратична похибка [Електронний ресурс] – режим доступу: [https://www.probabilitycourse.com/chapter9/mean\\_squared\\_error.php](https://www.probabilitycourse.com/chapter9/mean_squared_error.php)
18. Chollet Francois, "Deep Learning with Python", 2018. — 111 с.
19. Dense layer [Електронний ресурс] – режим доступу: [https://keras.io/api/layers/core\\_layers/dense/](https://keras.io/api/layers/core_layers/dense/)
20. BatchNormalization layer [Електронний ресурс] – режим доступу: [https://keras.io/api/layers/normalization\\_layers/batch\\_normalization/](https://keras.io/api/layers/normalization_layers/batch_normalization/)
21. Activation layer [Електронний ресурс] – режим доступу: <https://keras.io/api/layers/activations/>
22. Dropout layer [Електронний ресурс] – режим доступу: [https://keras.io/api/layers/regularization\\_layers/dropout/](https://keras.io/api/layers/regularization_layers/dropout/)
23. "Порівняння архітектурних стилів API" [Електронний ресурс] – режим доступу: <https://www.altexsoft.com/blog/soap-vs-rest-vs-graphql-vs-rpc/>
24. Офіційна документація FastAPI [Електронний ресурс] – режим доступу: <https://fastapi.tiangolo.com/>
25. Офіційна документація Flutter [Електронний ресурс] – режим доступу: <https://docs.flutter.dev/>
26. Офіційна сторінка PyCharm [Електронний ресурс] – режим доступу: <https://www.jetbrains.com/pycharm/>

27. Офіційна сторінка IntelliJ IDEA [Електронний ресурс] – режим доступу: <https://www.jetbrains.com/idea/>

28. Методичні вказівки до виконання економічної частини магістерських кваліфікаційних робіт / Уклад. : В. О. Козловський, О. Й. Лесько, В. В. Кавецький. – Вінниця : ВНТУ, 2021. – 42 с.

29. Кавецький В. В. Економічне обґрунтування інноваційних рішень: практикум / В. В. Кавецький, В. О. Козловський, І. В. Причепа – Вінниця : ВНТУ, 2016. – 113 с.

## **ДОДАТКИ**

## Додаток А (обов'язковий)

## Результат перевірки на плагіат в онлайн-системі UNICHECK



Ім'я користувача:  
Озеранський В.С. КН

ID перевірки:  
1013327397

Дата перевірки:  
19.12.2022 11:23:33 EET

Тип перевірки:  
Doc vs Internet + Library

Дата звіту:  
19.12.2022 11:28:16 EET

ID користувача:  
62038

Назва документа: 122МКР-БорисюкВМ2022

Кількість сторінок: 62 Кількість слів: 6872 Кількість символів: 52619 Розмір файлу: 1.55 MB ID файлу: 1013086949

Виявлено модифікації тексту (можуть впливати на відсоток схожості)

**11.6%**  
**Схожість**

Найбільша схожість: 5.97% з джерелом з Бібліотеки (ID файлу: 1013080093)

Не знайдено джерел з Інтернету

11.6% Джерела з Бібліотеки

3

Сторінка 64

**0% Цитат**

Не знайдено жодних цитат

Не знайдено жодних посилань

**1.41%**  
**Вилучень**

Деякі джерела вилучено автоматично (фільтри вилучення: кількість знайдених слів є меншою за 8 слів та 5%)

0.54% Вилучення з Інтернету

61

Сторінка 65

1.03% Вилученого тексту з Бібліотеки

164

Сторінка 65

**Модифікації**

Виявлено модифікації тексту. Детальна інформація доступна в онлайн-звіті.

Підозріле форматування

17  
сторінок

Рисунок А.1 - Результат перевірки на плагіат в онлайн-системі UNICHECK



## Додаток Б (обов'язковий)

### Лістинг програми

```

from sklearn.ensemble import RandomForestRegressor

RF_model = RandomForestRegressor(random_state=0, n_estimators=500,
max_depth=10)
RF_model.fit(X_train, y_train)

y_hat = RF_model.predict(X_test)
print ("R2 score:", r2_score(y_hat, y_test))
importances =
pd.DataFrame({'feature':X_train.columns,'importance':np.round(RF_model.feature
_importances_,3)})
importances =
importances.sort_values('importance',ascending=False).set_index('feature')
importances.plot.bar();
X_train, X_valid, y_train, y_valid = train_test_split(X, y, test_size=0.2)
params = {'num_leaves': 30,
          'min_data_in_leaf': 20,
          'objective': 'regression',
          'max_depth': 5,
          'learning_rate': 0.01,
          "boosting": "gbdt",
          "feature_fraction": 0.9,
          "bagging_freq": 1,
          "bagging_fraction": 0.9,
          "bagging_seed": 11,
          "metric": 'rmse',
          "lambda_11": 0.2,
          "verbosity": -1}
lgb_model = lgb.LGBMRegressor(**params, n_estimators = 10000, nthread = 4,
n_jobs = -1)
lgb_model.fit(X_train, y_train,
              eval_set=[(X_train, y_train), (X_valid, y_valid)], eval_metric='rmse',
              verbose=1000, early_stopping_rounds=200)
from sklearn import ensemble
params = {'n_estimators': 500, 'max_depth': 4, 'min_samples_split': 2,
          'learning_rate': .01, 'loss': 'ls'}
clf = ensemble.GradientBoostingRegressor(**params)
predictions2 = clf.fit(X_train,y_train)
training_score = clf.score(X_train, y_train)
print(f"Training Score: {training_score}")

```

```

predictions2 = np.expand_dims(clf.predict(X_test), axis = 1)
MSE = mean_squared_error(y_test, predictions2)
r2 = clf.score(X_test, y_test)
print(f"MSE: {MSE}, R2: {r2}")

import datetime as dt

from sqlmodel import SQLModel, Field, ForeignKey

class Movie(SQLModel, table=True):
    id: int = Field(primary_key=True)
    budget: int
    genre_ids: list[ForeignKey('genre.id')] = []
    homepage: str | None
    imdb_id: str
    original_language: str
    original_title: str
    overview: str
    popularity: float
    production_companies: list[ForeignKey('production_company.id')] = []
    production_countries: list[ForeignKey('production_country.id')] = []
    release_date: dt.datetime
    runtime: int
    keywords: list[ForeignKey('keyword.id')] = []
    cast: list[ForeignKey('actor.id')] = []
    crew: list[ForeignKey('credit.id')] = []
    revenue: int | None

from abc import abstractmethod, ABC
from typing import Iterable, Generic, TypeVar, NoReturn

from sqlmodel import select, delete, Session, SQLModel

from .models import Movie

_Model = TypeVar('_Model', bound=SQLModel)

class Repo(Generic[_Model], ABC):

    def __init__(self, session: Session):
        self._session = session

```

```
@abstractmethod
def get(self, pk: int) -> _Model:
    pass
```

```
@abstractmethod
def get_many(self, *where) -> Iterable[_Model]:
    pass
```

```
@abstractmethod
def delete(self, pk: int) -> NoReturn:
    pass
```

```
class MovieRepo(Repo[Movie]):
```

```
    def get(self, pk: int) -> Movie:
        query = select(Movie).where(Movie.id == pk)
        movie = self._session.execute(query).first()
        return movie
```


```
    def get_many(self, *where) -> Iterable[Movie]:
        query = select(Movie).where(*where)
        return self._session.execute(query)
```

```
    def delete(self, pk: int) -> NoReturn:
        query = delete(Movie).where(Movie.id == pk)
        self._session.execute(query)
```

## ІЛЮСТРАТИВНА ЧАСТИНА

ІНФОРМАЦІЙНА ТЕХНОЛОГІЯ РОЗПІЗНАВАННЯ  
МАТЕМАТИЧНИХ ФОРМУЛ НА ОСНОВІ НЕЙРОННОЇ  
МЕРЕЖІ

Виконав: студент 2-го курсу,  
групи КН-21м  
спеціальності 122 «Комп'ютерні науки»  
(шифр і назва напрямку підготовки, спеціальності)

 Борисюк В. М.

(прізвище та ініціали)

Керівник: к.т.н., доцент каф. КН

 Барабан С. В.

(прізвище та ініціали)

« 15 » 12 2022 р.

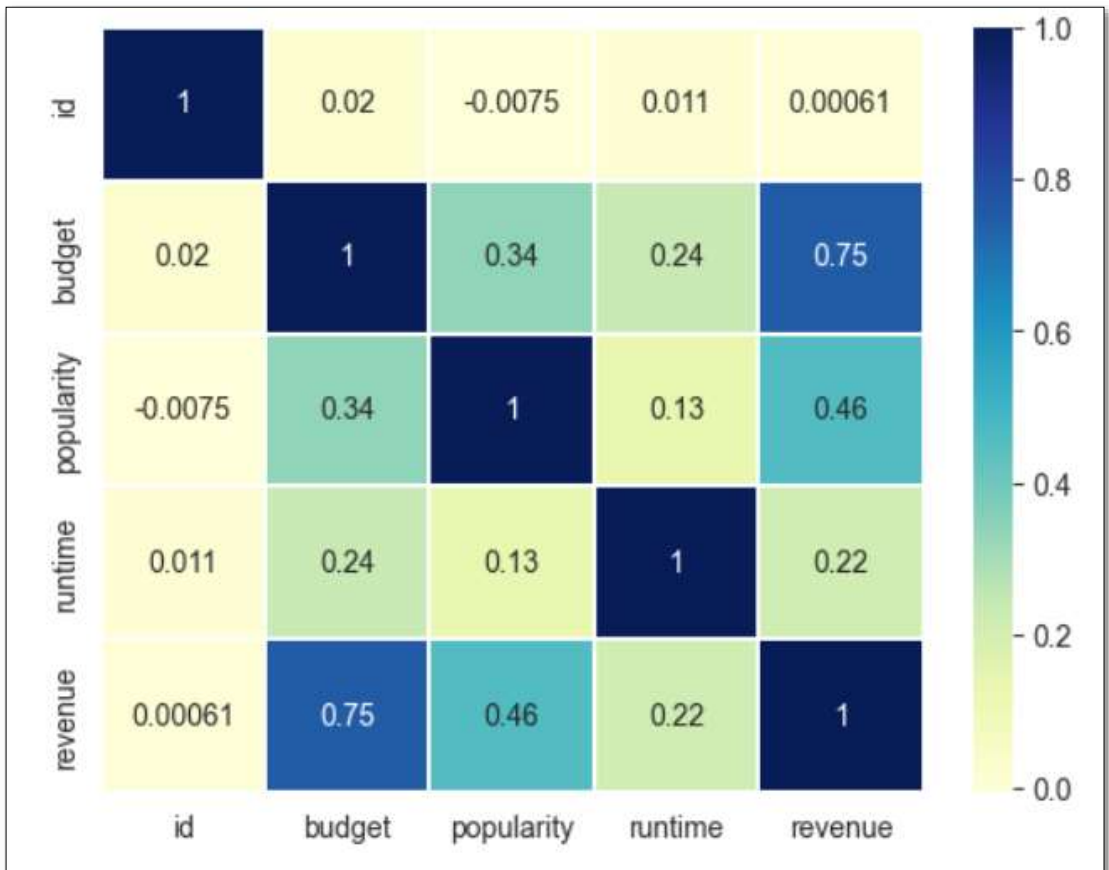


Рисунок В.1 – Графік кореляції ознак кінофільму до його доходу

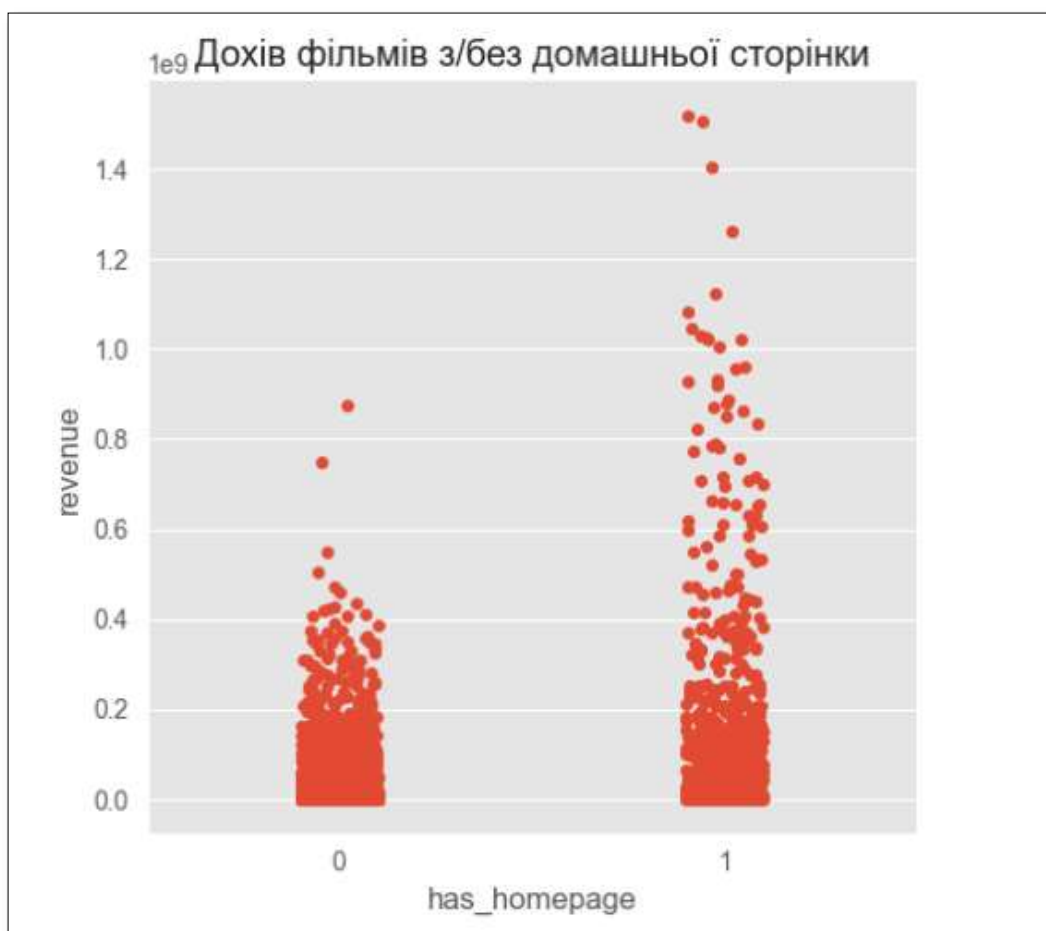


Рисунок В.2 – Графік розподілу доходу кінофільмів з домашньою і без неї

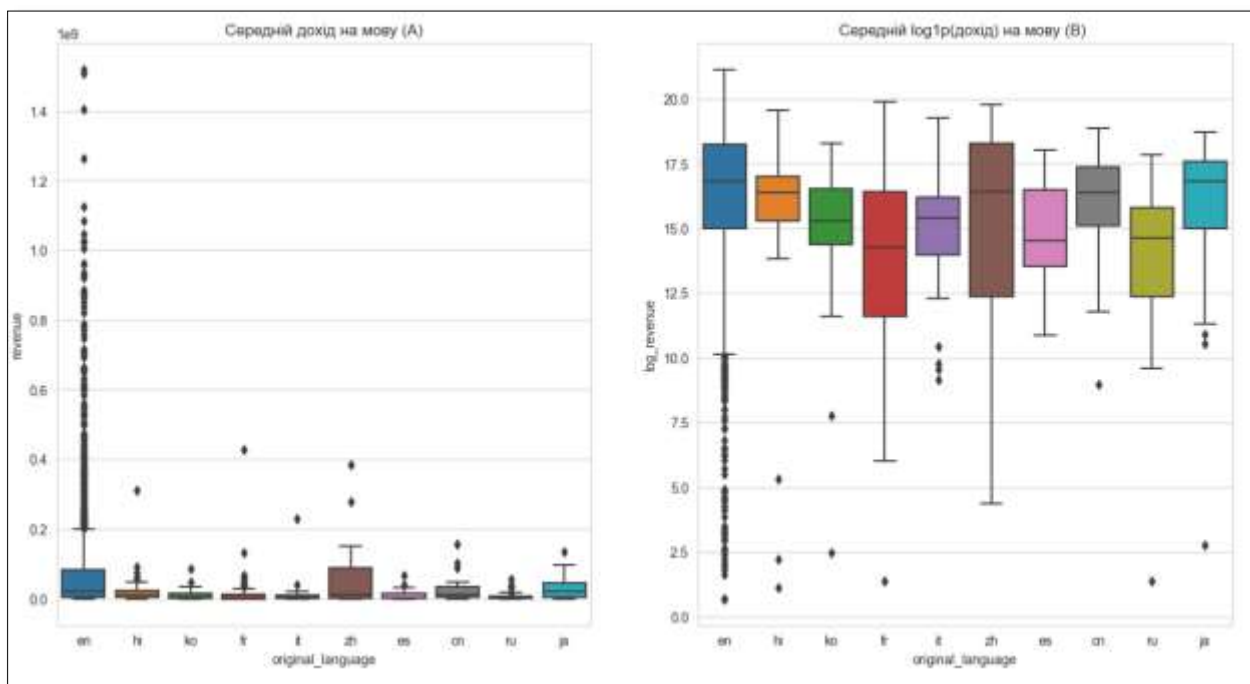


Рисунок В.3 – Графіки розподілу доходів для фільмів на різних мовах

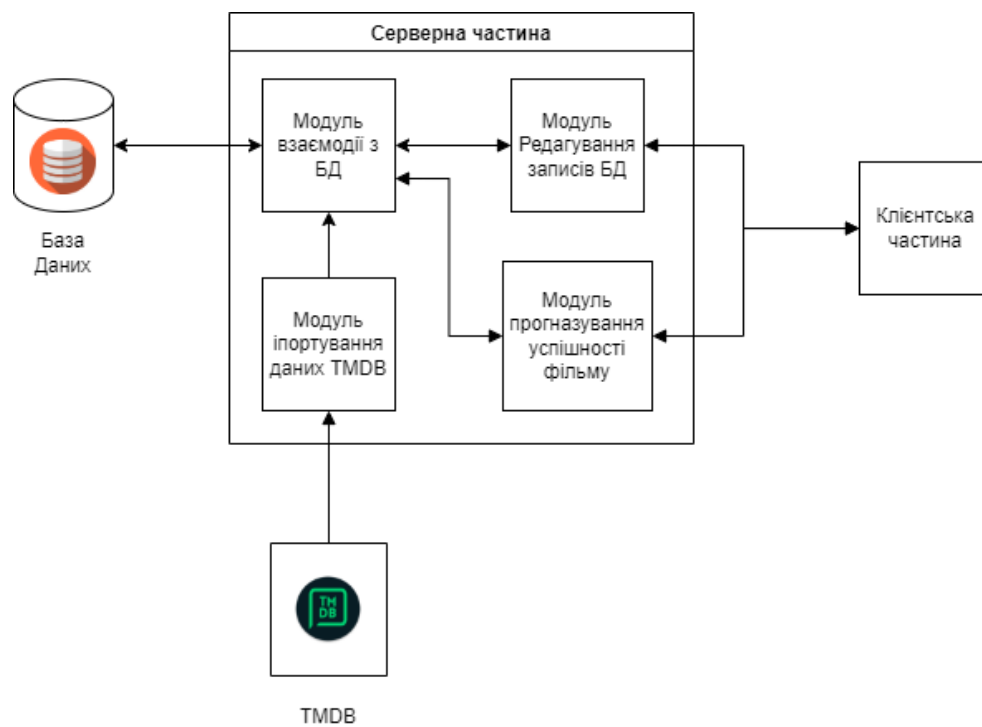


Рисунок В.4 – Загальна структура програмного забезпечення прогнозування успішності фільмів



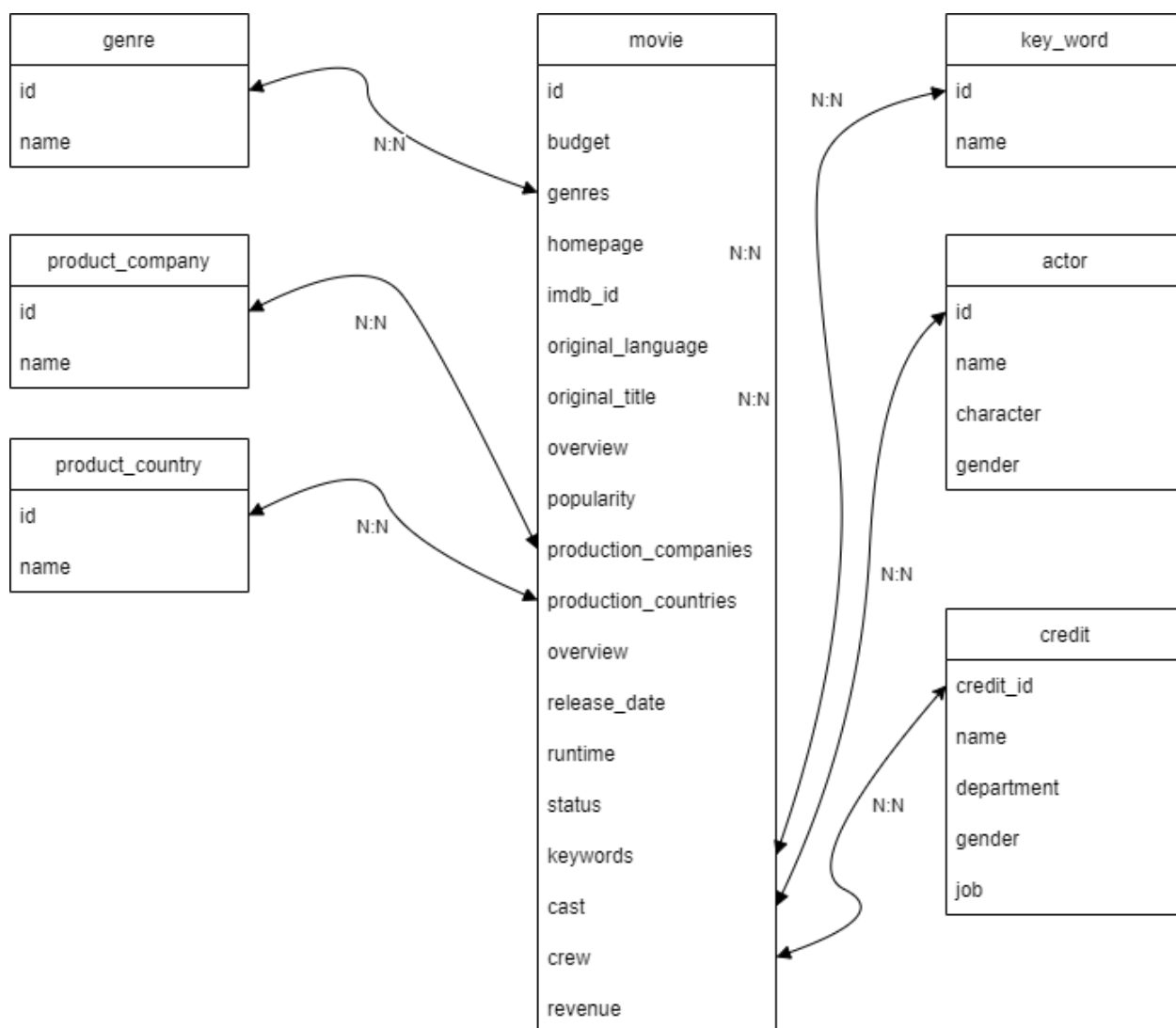


Рисунок В.5 - Загальна схема бази-даних кінофільмів

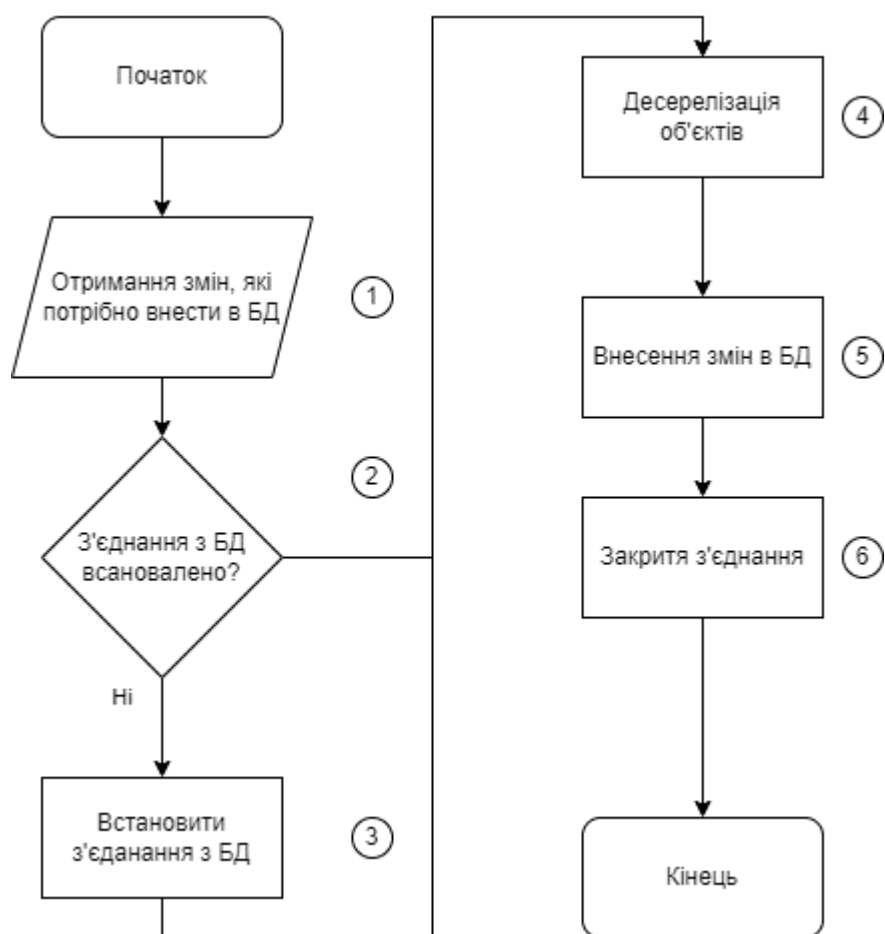


Рисунок В.6 - Схема алгоритму модуля взаємодії з базою даних

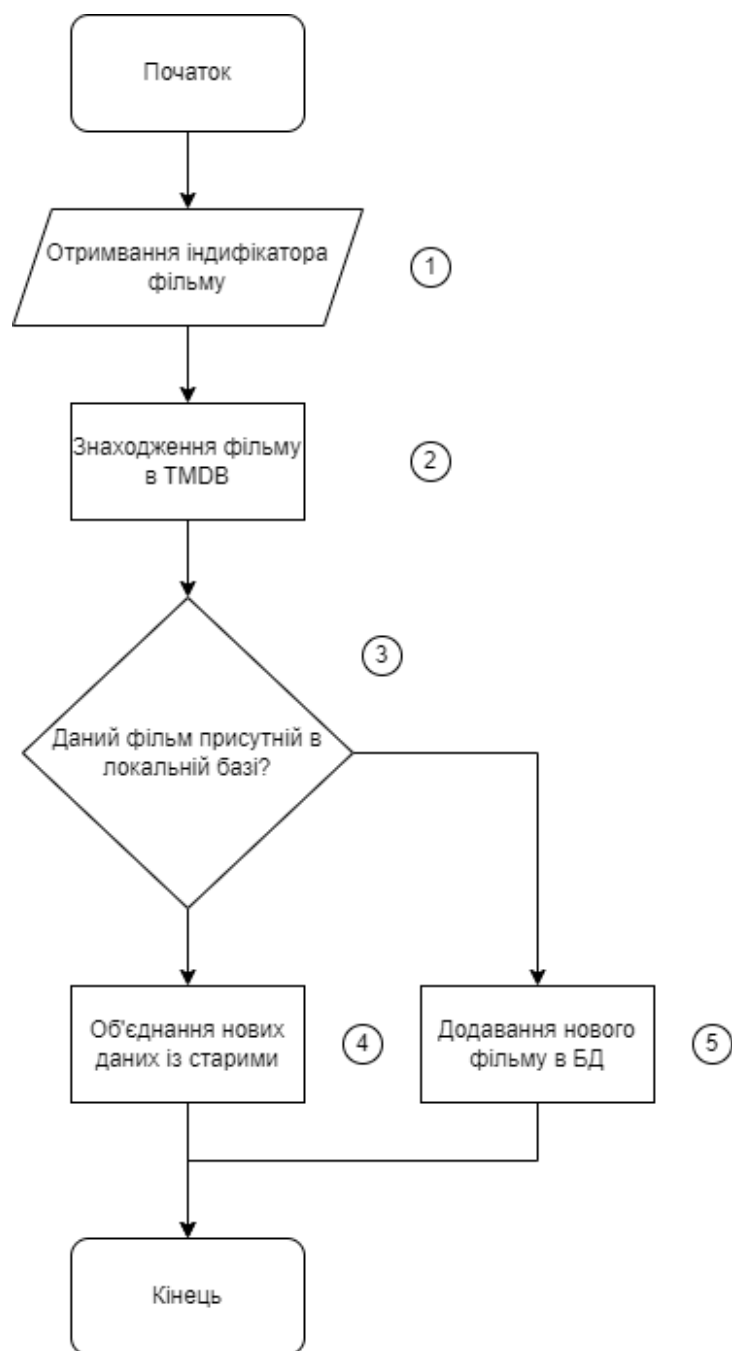



Рисунок В.7 – Схема алгоритму модуля імпортування фільмів з TMDB



Рисунок В.8 – Схема алгоритму модуля прогнозування успішності кінофільму

← The Batman DEBUG



Title: The Batman  
Director: Matt Reeves  
Budget: 200M \$

Revenue: 770M \$

Predicted revenue: 750M \$

Accuracy:

97 %

Рисунок В.9 – Вікно програми з результатом прогнозування успішності кінофільму

## Додаток Г (довідниковий)

### Інструкція користувача

Для використання програмного забезпечення, прогнозування успішності кінофільму, користувачу потрібно виконати наступні кроки:

1. Запустити ярлик програми – рисунок Г.1.

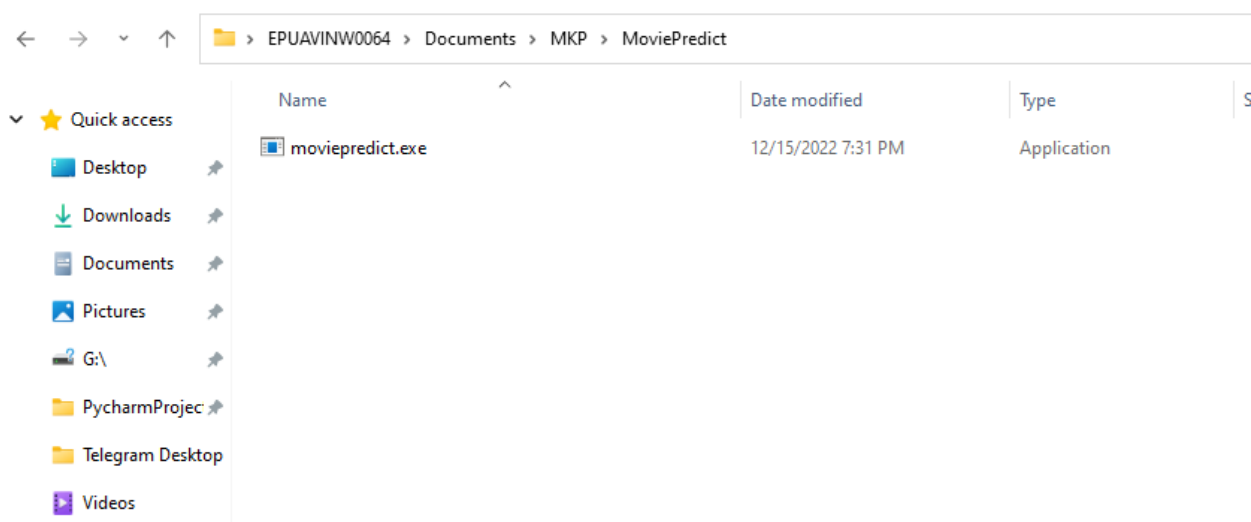


Рисунок Г.1 – Зображення ярлика програми

2. Відкриється головне вікно програмного забезпечення, на якому користувач може побачити увесь список кінофільмів- - рисунок Г.2.

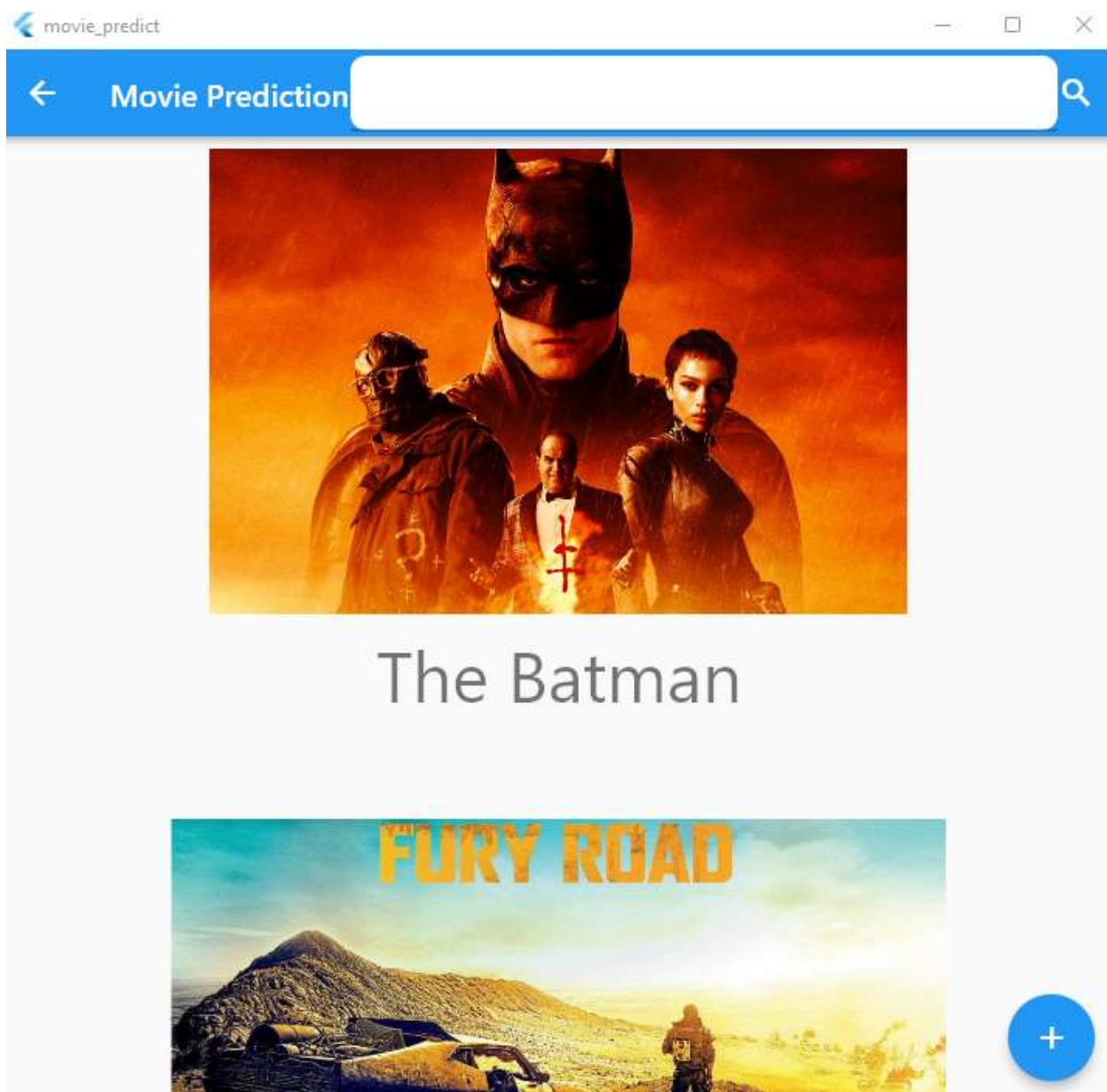


Рисунок Г.2 – Головне вікно програми прогнозування успішності кінофільму

3. Для пошуку фільмів в середині програми користувач може скористатись пошуковим полем у верхній частині вікна – рисунок Г.3.

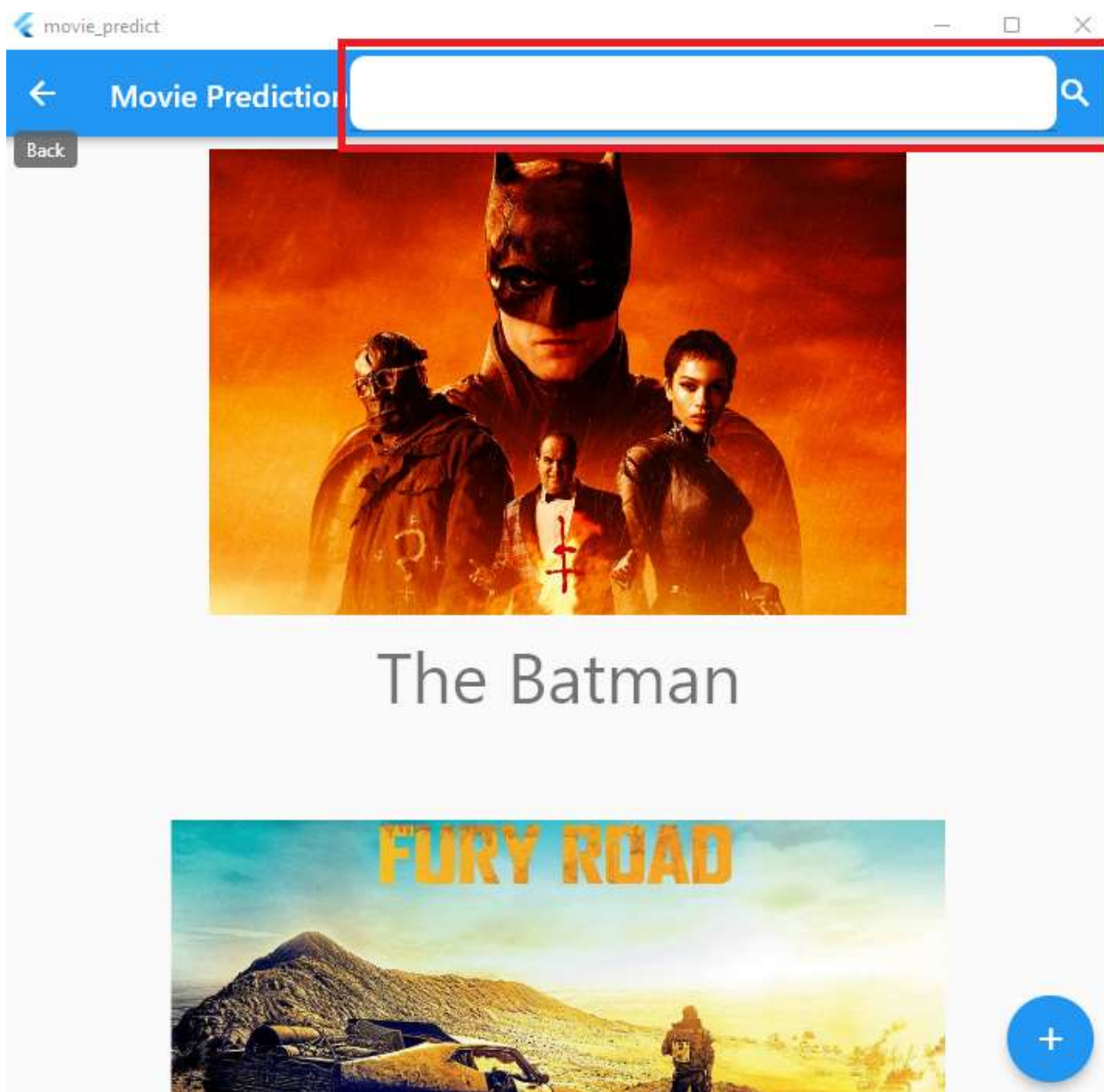



Рисунок Г.3 – Пошукове поле вводу

4. Щоб відкрити вікно з детальною інформацією про кінофільм та отримати прогноз щодо доходу кінофільму, потрібно натиснути на постер фільму. Причому, якщо інформація про реальний прибуток кінофільму присутня, виведеться індикатор, який відобразить точність прогнозу. Вікно з прогнозом доходу кінофільму на його детальною інформацією зображено на рисунку Г.4.



← The Batman



UNMASK THE TRUTH

Title: The Batman  
Director: Matt Reeves  
Budget: 200M \$

Revenue: 770M \$

Predicted revenue: 750M \$

Accuracy:

97 %

Рисунок Г.4 – Вікно з прогнозом доходу та детальною інформацією для кінофільму “The Batman”

5. Щоб додати новий кінофільм, потрібно натиснути відповідну кнопку на головному вікні – рисунок Г.5. Тоді відкриється вікно з формою, за допомогою якої можна додати новий кінофільм – рисунок Г.6.

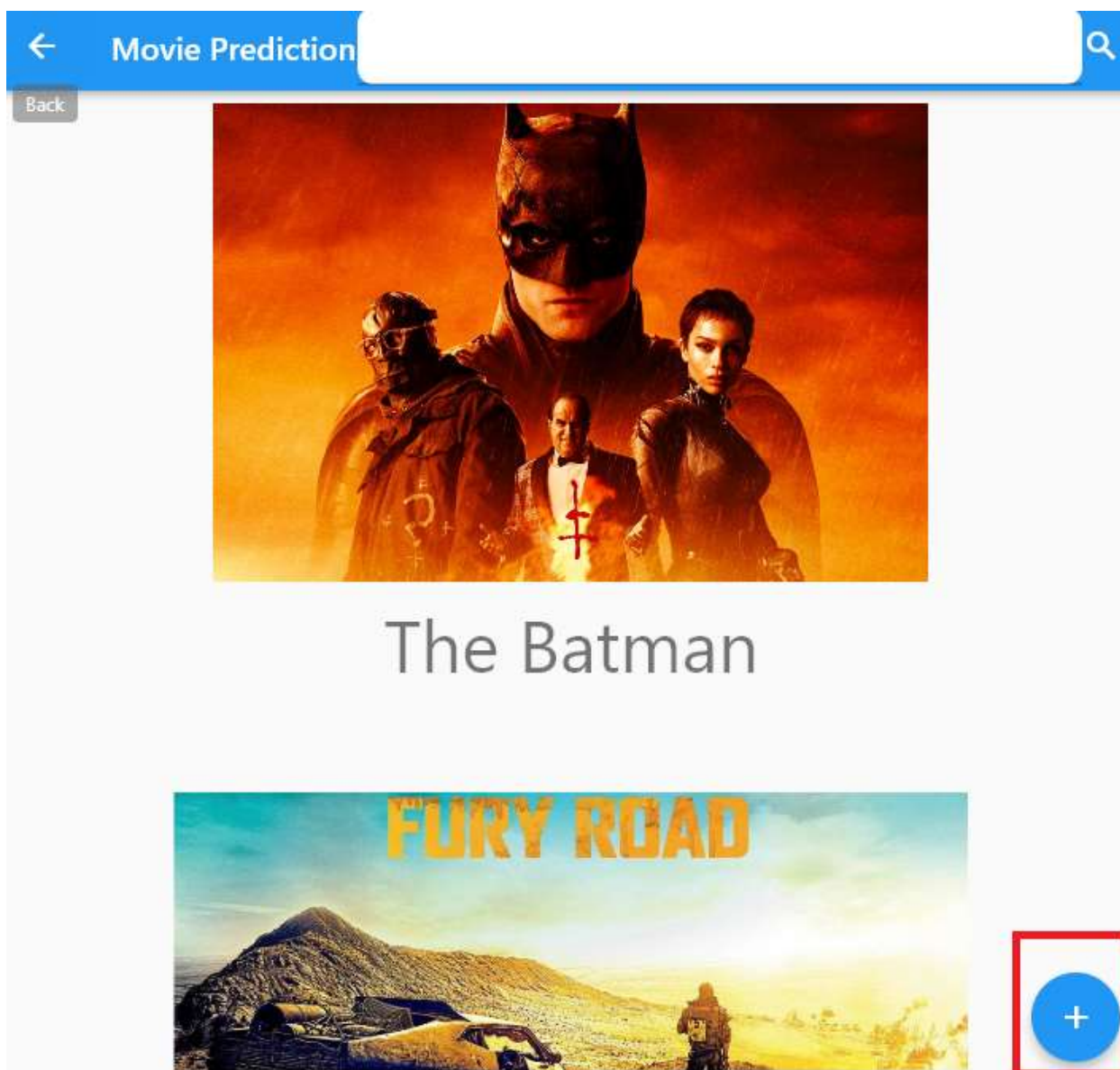
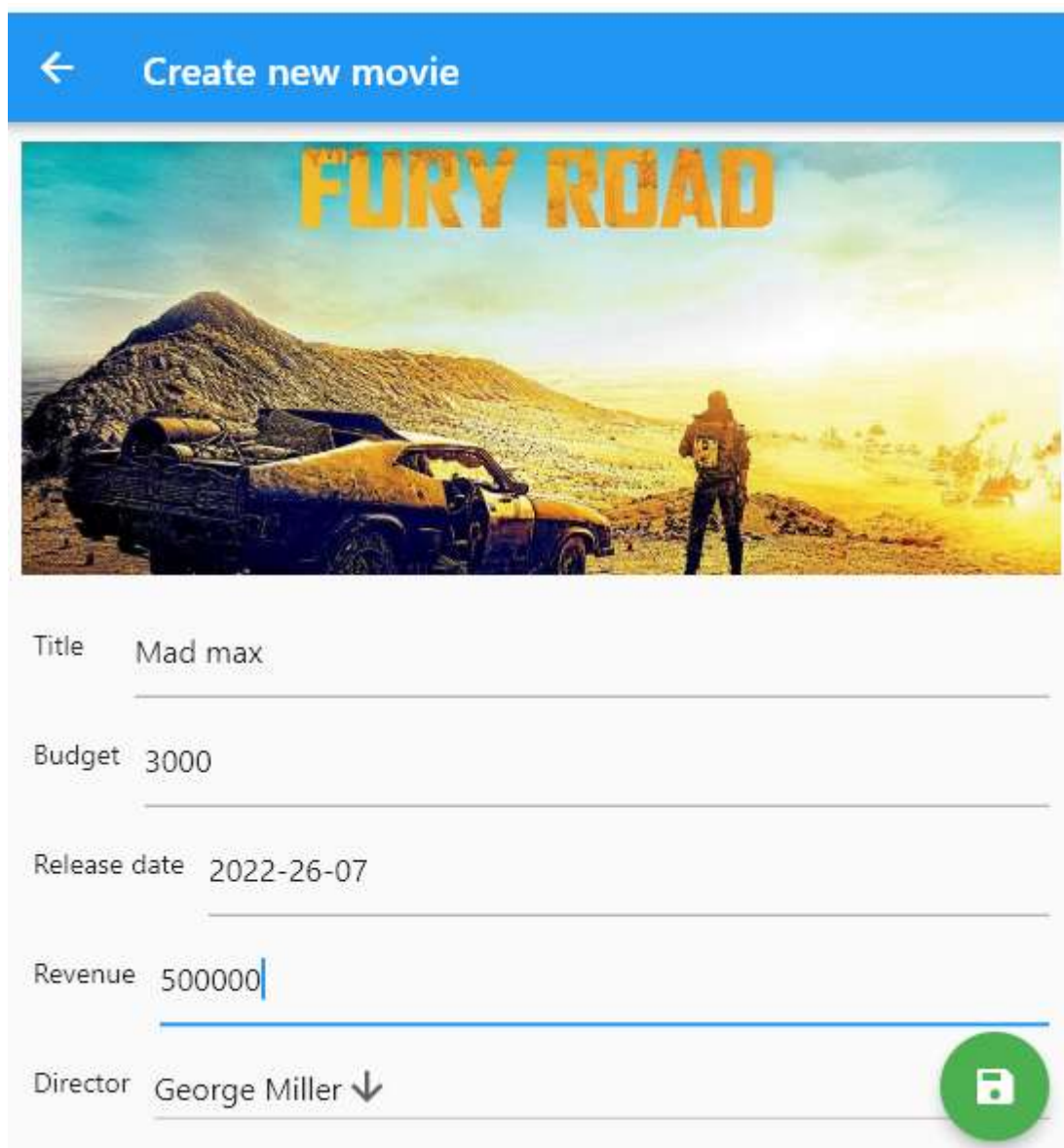


Рисунок Г.4 – Кнопка для переходу на форму для додавання нового кінофільму



← Create new movie

**FURY ROAD**

Title Mad max

Budget 3000

Release date 2022-26-07


Revenue 500000

Director George Miller ↓

Рисунок Г.5 – Форма для додавання нового кінофільму

6. Для збереження кінофільму потрібно натиснути відповідно кнопку в правому нижньому куті – рисунок Г.6.

← Create new movie



Title Mad max

Budget 3000

Release date 2022-26-07

Revenue 500000

Director George Miller ↓




Рисунок Г.6 – Кнопка збереження кінофільму