

МІНІСТЕРСТВО ОСВІТИ І НАУКИ УКРАЇНИ
Західноукраїнський національний університет
Факультет комп'ютерних інформаційних технологій
Кафедра інформаційно-обчислювальних систем і управління

МИСЬ Андрій Русланович

**Інтелектуальний модуль прогнозування врожайності ягід /
Intelligent berry yield prediction module**

Спеціальність 122 – Комп'ютерні науки
Освітньо-професійна програма – Комп'ютерні науки

Дипломний проект

Виконав студент групи КН-41
А. Р. Мись

Науковий керівник:
к.т.н., доцент Т.В. Лендюк

Дипломний проект допущено до захисту
«___» _____ 2023 р.

Завідувач кафедри
_____ М.П. Комар

Тернопіль – 2023

Факультет комп'ютерних інформаційних технологій
Кафедра інформаційно-обчислювальних систем і управління
Освітній ступінь «бакалавр»
Спеціальність 122 – Комп'ютерні науки
Освітньо-професійна програма – Комп'ютерні науки

ЗАТВЕРДЖУЮ
Завідувач кафедри
_____ М.П. Комар
« _____ » _____ 2022р.

З А В Д А Н Н Я
НА ДИПЛОМНИЙ ПРОЕКТ СТУДЕНТУ
Мисю Андрію Руслановичу
(прізвище, ім'я, по батькові)

1. Тема проекту: Інтелектуальний модуль прогнозування врожайності ягід /
Intelligent berry yield prediction module
керівник проекту к.т.н., доцент Т.В. Лендюк

затверджені наказом по університету від 08 грудня 2022 р. № 491.

2. Строк подання студентом закінченого проекту 01 червня 2023 р.

3. Вихідні дані до проекту: технічне завдання.

4. Зміст розрахунково-пояснювальної записки (перелік питань, які потрібно розробити)

- Вивчення та аналіз сфери дослідження
- Перегляд та аналіз існуючих рішень в даній області
- Оцінка методології ансамблевого навчання
- Розгляд XGBoost, LightGBM та CatBoost як потенційних інструментів для прогнозування врожайності ягід
- Розробка алгоритму навчання моделей для прогнозування врожайності ягід
- Розробка алгоритму валідації моделей для перевірки точності прогнозування
- Реалізація та валідація ансамблевого навчання моделей для прогнозування врожайності ягід

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

- алгоритм навчання моделей для прогнозування врожайності ягід
- алгоритм валідації моделей для перевірки точності прогнозування

6. Консультанти розділів проекту

Розділ	Прізвище, ініціали та посада консультанта	Підпис, дата	
		Завдання видав	Завдання прийняв
Н. контроль	к.т.н., доцент Т.В. Лендюк		

7. Дата видачі завдання 08 грудня 2022 р.

КАЛЕНДАРНИЙ ПЛАН

№ з/п	Назва етапів дипломного проекту	Строк виконання етапів проекту	Примітка
1	Аналіз предметної області і постановка задачі дослідження	30.12.2022	
2	Алгоритмічне забезпечення інтелектуального модуля прогнозування врожайності ягід	24.03.2023	
3	Програмно-технологічне забезпечення інтелектуального модуля прогнозування врожайності ягід	12.05.2023	
4	Повне завершення та оформлення дипломного проекту	01.06.2023	

Студент _____ А. Р. Мись
(підпис)

Керівник проекту _____ Т.В. Лендюк
(підпис)

РЕФЕРАТ

Пояснювальна записка до дипломного проекту: 62 с., 7 рис., 1 табл., 2 додатки, 48 джерел.

Метою дипломного проекту є розробка інтелектуального модуля прогнозування врожайності ягід, заснованого на алгоритмах машинного навчання.

Об'єктом дослідження є процес прогнозування врожайності ягід.

Предметом дослідження є використання ансамблевих моделей машинного навчання, зокрема XGBoost, LightGBM та CatBoost, для створення точного прогнозу врожайності ягід.

Розроблено та досліджено програмне забезпечення для інтелектуального модуля прогнозування врожайності ягід.

Розроблений модуль може допомогти фермерам і аграріям планувати свою діяльність, оптимізувати витрати та підвищувати продуктивність. Модуль також може стати корисним інструментом для дослідників і аналітиків в агросекторі.

ПРОГНОЗУВАННЯ ВРОЖАЙНОСТІ ЯГІД, АЛГОРИТМИ МАШИНОГО НАВЧАННЯ, АНСАМБЛЕВЕ НАВЧАННЯ, XGBOOST, LIGHTGBM, CATBOOST, ВАЛІДАЦІЯ МОДЕЛЕЙ.

ABSTRACT

The bachelor's thesis report: 62 pages, 7 figures, 1 table, 2 appendices, 48 references.

The aim of the diploma project is to develop an intelligent module for predicting the yield of berries based on machine learning algorithms.

The object of the research is the process of predicting the yield of berries.

The subject of the research is the use of ensemble machine learning models, specifically XGBoost, LightGBM, and CatBoost, to create an accurate forecast of berry yield.

Software has been developed and studied for the intelligent module for predicting the yield of berries.

The developed module can assist farmers and agriculturalists in planning their activities, optimizing costs, and increasing productivity. The module can also be a useful tool for researchers and analysts in the agricultural sector.

BERRY YIELD PREDICTION, MACHINE LEARNING ALGORITHMS, ENSEMBLE LEARNING, XGBOOST, LIGHTGBM, CATBOOST, MODEL VALIDATION.

ТЕХНІЧНЕ ЗАВДАННЯ

1. НАЙМЕНУВАННЯ ТА ОБЛАСТЬ ЗАСТОСУВАННЯ

1.1 Інтелектуальний модуль прогнозування врожайності ягід.

1.2 Область застосування – сільське господарство та агросектор.

2. ОСНОВА ДЛЯ РОЗРОБЛЕННЯ

Основою для розроблення є завдання на дипломний проект, затверджене кафедрою інформаційно-обчислювальних систем і управління факультету комп'ютерних інформаційних технологій Західноукраїнського національного університету.

3. ПРИЗНАЧЕННЯ РОЗРОБЛЕНОГО КОМПЛЕКСУ

Метою дипломного проекту є розробка інтелектуального модуля прогнозування врожайності ягід, заснованого на алгоритмах машинного навчання. Цей модуль покликаний підвищити точність прогнозування врожайності, враховуючи велику кількість взаємопов'язаних факторів.

4. ДЖЕРЕЛА РОЗРОБЛЕННЯ

Джерелами даної розробки є матеріали навчальної і реферативної літератури, технічна документація, науково-дослідні статті, журнали, Інтернет.

5. ТЕХНІЧНІ ВИМОГИ

5.1 Основні функціональні вимоги до програмної системи:

– вибір даних (пакетних даних / даних у реальному часі) з різних архітектур великих даних;

– провести попередню обробку даних для інтелектуального аналізу тексту.

– Реалізувати та провести валідацію ансамблевого навчання моделей для прогнозування врожайності ягід.

5.2 Вимоги до апаратних засобів:

– кластер запускається на фізичній машині Macbook Pro з 16 гігабайт оперативної пам'яті та процесором 2.3 GHz Intel Core i5.

5.3 Вимоги до програмних засобів:

- для розробки програмне забезпечення - Python 3.7;
- для створення графічного інтерфейсу користувача використано – tkinter, Adobe XD;
- для реалізації моделей навчання – фреймворк XGBoost, LightGBM та CatBoost.

6. ПОРЯДОК КОНТРОЛЮ

6.1 Представлення дипломного проекту на попередній захист.

6.2 Представлення дипломного проекту на захист.

Завдання прийняв до виконання _____ А. Р. Мись
(підпис) (прізвище та ініціали)

Керівник дипломного проекту _____ Т.В. Лендюк
(підпис) (прізвище та ініціали)

ЗМІСТ

Вступ.....	9
1 Аналіз предметної області і постановка задачі дослідження	11
1.1 Аналіз предметної області	11
1.2 Огляд існуючих рішень	14
1.3 Оцінка підходів ансамблевого навчання	16
1.4 Постановка задачі дослідження.....	24
2 Алгоритмічне забезпечення інтелектуального модуля прогнозування врожайності ягід.....	26
2.1 Огляд XGBoost, LightGBM та CatBoost як інструментів для прогнозування врожайності ягід	26
2.2 Алгоритм навчання моделей прогнозування врожайності ягід...	29
2.3 Алгоритм валідації моделей прогнозування врожайності ягід....	31
3 Програмно-технологічне забезпечення інтелектуального модуля прогнозування врожайності ягід	35
3.1 Розвідувальний аналіз набору даних	35
3.2 Ансамблеве навчання моделей.....	40
3.3 Валідація моделей.....	48
Висновки	55
Список використаних джерел.....	57
Додаток А.....	63
Додаток Б	64

					<i>ДП.КН.8351441.062.ПЗ</i>
<i>Змн.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дат</i>	
<i>Розроб.</i>		<i>Мись А. Р.</i>			<i>Літ.</i>
<i>Перевір.</i>		<i>Лендюк Т.В.</i>			<i>Аркуш</i>
<i>Консульт.</i>					8
<i>Н. Контр.</i>		<i>Лендюк Т.В.</i>			<i>Аркушів</i>
<i>Затверд.</i>		<i>Комар М.П.</i>			62
					<i>ЗУНУ.ФКІТ.КН-41</i>

ВСТУП

Актуальність теми. Прогнозування врожайності ягід є важливим аспектом сільськогосподарської промисловості. Це не тільки допомагає фермерам у плануванні та оптимізації виробництва, але також сприяє стабільності ринку, забезпечуючи достатній обсяг продукції для задоволення попиту.

Однак, врожайність ягід може відчутно відрізнятись від року до року через різноманітність факторів, включаючи погодні умови, шкідники та хвороби, а також практики управління. Тому розробка надійних методів прогнозування врожайності стає все більш актуальною.

Використання технологій штучного інтелекту, таких як алгоритми машинного навчання, може значно покращити точність прогнозів, враховуючи велику кількість факторів та їх взаємодію.

Це особливо актуально в умовах глобальних змін клімату, які можуть призвести до непередбачуваних відхилень у врожайності ягід. За допомогою штучного інтелекту, фермери та сільськогосподарські підприємства можуть адаптуватися до цих змін, оптимізуючи свої стратегії для забезпечення максимальної врожайності.

Таким чином, розробка інтелектуального модуля прогнозування врожайності ягід є актуальною та важливою задачею, яка може мати значний вплив на сільськогосподарську промисловість та продовольчу безпеку.

Метою дипломного проекту є розробка інтелектуального модуля прогнозування врожайності ягід, заснованого на алгоритмах машинного навчання. Цей модуль покликаний підвищити точність прогнозування врожайності, враховуючи велику кількість взаємопов'язаних факторів.

Для досягнення поставленої мети було необхідно провести ряд теоретичних досліджень, визначити та послідовно вирішити наступні завдання:

1. Вивчення та аналіз сфери дослідження

									ДП.КН. 8351441.062ПЗ	Арк.
										9
Зм.	Арк.	№ докум.	Підпис	Дата						

1 АНАЛІЗ ПРЕДМЕТНОЇ ОБЛАСТІ І ПОСТАНОВКА ЗАДАЧІ ДОСЛІДЖЕННЯ

1.1 Аналіз предметної області

Метою інтелектуального модуля прогнозування врожайності ягід є аналіз та прогнозування врожайності лохини. Лохина є плодовою рослиною, яка вирощується для отримання ягід, які мають велике значення в харчовій та фармацевтичній промисловості.

Дика лохина є найважливішою плодовою культурою штату Мен, яка росте на гірських кислих піщаних ґрунтах. Дика лохина (також відома як низькокущова лохина) є найважливішою культурою для більшості людей і вважається найбільшим виробником серед інших культур у США [1]. Головна перевага дикої лохини перед культурною полягає в тому, що вона не садиться, а росте природним шляхом на кам'янистих гірках і лісах. Культурну лохину, навпаки, необхідно садити [2]. Дика чорниця дає різноманітні культури, оскільки вона захищена від селекційної селекції та генної інженерії. Ручні комбайни використовують граблі, щоб зішкребти ягоди з рослин [3] через пересічену місцевість і висоту дикоростучої чорниці. Найкращий час для збору дикої чорниці в штаті Мен — з кінця липня до початку вересня. [4]. Види комах-запилювачів і погодні умови є основними факторами, що впливають на врожайність дикої лохини [5]. Кількість комах-запилювачів, оптимальний час для запилення дикої лохини (близько 6:00 – 17:00), види бджіл (медоносна бджола (НВ), джміль (ВВ), Андрена (АД) і *Osmia* (OS)) та погодні умови. (низькі температури, волога погода, сильний вітер, заморозки та підвищення температури) є факторами, що впливають на врожайність чорниці дикої [6]. Отримати надійні прогнози врожайності для досягнення цілей Продовольчої та сільськогосподарської організації (ФАО) важко. Основною метою ФАО є викорінення глобального голоду. Кілька

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		11

вчених зосередилися на точному землеробстві шляхом прогнозування врожайності за допомогою методології машинного навчання (ML) [7] – [10]. Однак лише кілька досліджень використовували EMLA з різними топологіями [11], [12]. Для MLA потрібні великі та високоякісні набори даних для точної оцінки врожайності [13]. Моделювання систем сільськогосподарського виробництва має важливе значення, особливо коли збір великомасштабних часових і географічних даних є надзвичайно дорогим, складним або в деяких випадках неможливим.

У кількох дослідженнях якість даних, які використовуються під час навчання, є недостатньою, що може вплинути на результати експерименту [14], що погіршує точність прогнозу врожайності [15]. Щоб вирішити ці проблеми, були використані підходи комп'ютерного імітаційного моделювання, які імітують вимірювання на основі сцен для створення даних [16]. Для підтвердження комп'ютерної імітаційної моделі тестуються припущення моделі. Аналіз чутливості використовується для припущень даних. Його використовували для демонстрації інтерактивного впливу домінування бджолиної спільноти та метеорологічних змінних на вихід дикої чорниці [17]. Центр моделювання, експериментів і ефективного дизайну обробки даних у Військово-морських аспірантурах визначає обробку даних як практику накопичення наборів даних за допомогою моделювання та комп'ютерного моделювання.

Таким чином, мета обробки даних полягає в тому, щоб дати особам, які приймають рішення, зрозуміти складні проблеми за допомогою моделювання для генерування даних [18].

Obsie та ін. [17] використовували обчислювальні експерименти для отримання навчальних даних, які згодом можуть бути використані для навчання та оцінки MLA.

Для емпіричних наборів даних EMLA перевершила індивідуальні MLA з точки зору точності. Основна перевага використання EMLA полягає в тому, що вони можуть вирішувати складні проблеми [19]. Основна ідея

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		12

використання EMLA полягає в тому, щоб перетворити слабкого учня на потужного для покращення підтримки прийняття рішень [20]. Кілька досліджень поєднували різні методи МН для досягнення високої точності. Методи стекування використовуються для реалізації моделей класифікації, таких як PredT4SE-Stack [21] і NeuroPpred-Fuse [22], а також зважених методів стекування кількох метамodelей [23] для регресії. Кілька методів з різними архітектурами використовувалися для допомоги точному землеробству через прогнозування врожайності [12]. Техніка стекування отримує всі вихідні дані передбачення з інших МЛА і передає їх у метамodelь для отримання найкращого рішення. Автори [17] довели, що дані моделювання в поєднанні з кліматичними даними корисні для навчання МЛА, оскільки ці дані мають високу якість і вони створюють набори даних, які важко зібрати вручну. Основна мета полягає в тому, щоб уникнути дорогих наборів даних у реальному часі, отриманих із супутникових зображень із середньою вартістю США \$ 1/км² [24].

Цей модуль використовує дані про різні фактори, що впливають на врожайність лохини, такі як розмір клону лохини, густина бджіл (медоносних та диких), температурні умови, кількість днів з опадами та інші. З використанням цих даних модуль проводить аналіз та прогнозує врожайність лохини.

Аналіз предметної області дозволяє зрозуміти, які фактори мають найбільший вплив на врожайність лохини і як ці фактори взаємодіють між собою. Це допомагає визначити оптимальні умови для вирощування лохини та розробити стратегії для підвищення врожайності.

Інтелектуальний модуль прогнозування врожайності ягід лохини є потужним інструментом для сільськогосподарських досліджень, розробки нових методів вирощування та оптимізації врожайності лохини.

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		13

1.2 Огляд існуючих рішень

У країнах, що розвиваються, урожайність зазвичай оцінюють вручну. Однак інші країни зосередилися на використанні нових технологій, включаючи MLA та алгоритми глибокого навчання, для прогнозування врожайності.

Протягом 30 років на західному узбережжі Індії збиралися часові ряди даних про врожайність рису в поєднанні з метеорологічними характеристиками. Чотири MLA були використані для прогнозування врожайності рису. Загалом, оператор найменшого абсолютного скорочення та відбору (LASSO) досяг найкращих результатів [25].

Гібридна множинна лінійна регресія (MLR) – штучна нейронна мережа (ANN) була використана для прогнозування врожайності [26]. У цій мережі, замість довільної ініціалізації ваги та зміщення, ваги вхідного рівня та зміщення ініціалізуються з використанням коефіцієнтів MLR та зміщення. Площа посіву, регіон водної системи, споживання добрив і точки інтересу водної системи були включені до створених сільськогосподарських даних. Серед кліматичних характеристик, включених до набору даних, були кліматичні опади, найвищі та найнижчі температури та сонячна радіація. MLR–ANN було застосовано до цих наборів даних для точного прогнозування врожайності рису.

П'ять MLA, а саме лінійна регресія (LR), LASSO, XGBoost, LGBM і випадковий ліс (RF), були об'єднані за допомогою оптимізованого зваженого ансамблю для прогнозування даних про врожайність кукурудзи, включаючи управління культурами, клімат, умови ґрунту та історичну врожайність. П'ять MLA у поєднанні з оптимізованим зваженим ансамблем перевершили існуючі MLA [12].

Шаххоссейні та ін. [27] представили нову ансамблеву згорову нейронну мережу (CNN). Крім того, глибокі нейронні мережі використовувалися для прогнозування врожайності кукурудзи. Було розглянуто набір даних, що

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		14

містить дату посіву, дані про погоду, характеристики ґрунту та історичну врожайність кукурудзи.

Урожайність зернових була передбачена за допомогою MLA, навченого з набором даних, що містить дані про погоду, дистанційне зондування, клімат і дані про врожайність [28], і набір даних був об'єднаний із супутниковими даними та історичною врожайністю. Ці чотири вимірювання параметрів були об'єднані, щоб сформувати єдиний набір даних. Було використано одну лінійну MLA і три нелінійні MLA. XGBoost перевершив інші MLA для цієї програми.

Лінійні та нелінійні моделі були використані для оцінки врожайності сої та кукурудзи за допомогою зображень Landsat та SPOT. Було вивчено дві моделі, нейронну мережу та MLR [29], і результати показали, що коли були включені всі зображення, моделі створили сильну відповідність, перевищуючи нейронні мережі з точки зору точності.

Хакі та ін. [30] представили YieldNet, яка є сучасною моделлю CNN, яка поєднує в собі передбачувану врожайність кукурудзи та сої шляхом перемикання ваги хребта, включаючи екстрактор, використовуючи новий метод глибокого навчання. Набір даних було створено за допомогою даних дистанційного зондування, що дозволяє безперервно контролювати культуру протягом усього циклу росту.

Набори даних були створені за допомогою імітаційної моделі запилення дикої чорниці. Ця інформація містить дані про чотири види бджіл, розміри їхніх клонів (CS) і погодні дані. Для прогнозування врожайності дикої чорниці використовували чотири MLA. XGBoost [17] досяг найкращої продуктивності за всіма показниками оцінювання.

В процесі аналізу було зроблено огляд декількох досліджень, в яких використовувалися дані моделювання для прогнозування врожайності. Зокрема, статті [17] та [31] досліджували врожайність кукурудзи протягом 5 років за допомогою симулятора систем сільськогосподарського виробництва в Середньому Заході США.

									ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата						15

Ці дослідження показали значну перевагу використання ансамблевого моделювання при прогнозуванні врожайності. Ансамблеві моделі комбінують прогнози кількох базових моделей, що дозволяє отримати більш точний та стабільний результат. Вони враховують різноманітність даних та різні підходи до моделювання, що сприяє зменшенню помилок та підвищенню точності прогнозування.

Ансамблеві моделі мають широке застосування в сільському господарстві та дозволяють вирішувати складні задачі, пов'язані з прогнозуванням врожайності. Використання таких моделей допомагає зрозуміти вплив різних факторів на врожайність та розробити ефективні стратегії вирощування.

Таким чином, ансамблеве моделювання є потужним інструментом у дослідженнях прогнозування врожайності, оскільки воно забезпечує більш точні та надійні результати, що можуть бути використані для покращення сільськогосподарського виробництва.

1.3 Оцінка підходів ансамблевого навчання

Ансамблювання є потужним підходом в машинному навчанні, який дозволяє поєднувати прогнози декількох моделей з метою отримання кращої і точнішої прогнозової здатності. Ідея полягає в тому, що комбінування прогнозів з різних моделей може знизити помилку та покращити загальну проціс прогнозування.

Основна мета ансамблювання полягає у використанні комбінації слабких моделей (тобто моделей з невеликою прогнозовою здатністю) для створення потужної моделі, яка може робити кращі прогнози. Це досягається шляхом використання різних стратегій комбінування, які можуть варіюватися залежно від методу ансамблювання.

Основні методи ансамблювання включають:

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		16

1. Бутстреп-агрегація (Bagging): цей метод полягає у використанні декількох незалежних моделей навчання на різних підмножинах тренувальних даних замість однієї моделі.

Переваги бутстреп-агрегації (Bagging):

- Зменшення дисперсії: Використання декількох незалежних моделей дозволяє зменшити варіацію прогнозів. Кожна модель навчається на різних підмножинах даних, тому вони можуть виявляти різні взаємозв'язки в даних.
- Уникнення перенавчання: Бутстреп-агрегація допомагає уникнути перенавчання, оскільки моделі навчаються на різних підмножинах даних. Це дозволяє їм більш загально узагальнювати залежності в даних.
- Врахування різноманітності: Кожна модель навчається на випадково вибраних підмножинах даних, що сприяє різноманітності моделей. Це може поліпшити прогностичні здібності ансамблю.

Недоліки бутстреп-агрегації:

- Збільшення обчислювальних витрат: Оскільки бутстреп-агрегація використовує декілька незалежних моделей, це може призвести до збільшення обчислювальних витрат, особливо якщо модель досить складна або навчається на великому обсязі даних.
- Затримка в навчанні: Тренування декількох моделей може зайняти більше часу, ніж тренування однієї моделі. Це може бути проблемою, коли швидкість тренування є критичним фактором.
- Втрата інтерпретованості: Бутстреп-агрегація використовує незалежні моделі, тому важче інтерпретувати результати прогнозів порівняно з однією моделлю. Це може бути проблемою, якщо потрібно зрозуміти вплив окремих змінних на прогнози.
- Врахування ваги: Бутстреп-агрегація може надати вагу кожній моделі в ансамблі, залежно від її ефективності. Це дозволяє приділити

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		17

більшу увагу до добре працюючих моделей і меншу увагу до менш ефективних моделей.

- Використання для класифікації та регресії: Бутстреп-агрегація може застосовуватись як для задач класифікації, так і для задач регресії. Це робить його універсальним методом для прогнозування в різних задачах.
- Відповідність до паралельного обчислення: Бутстреп-агрегація може бути ефективно розпаралелена, оскільки кожна модель може бути навчена незалежно. Це дозволяє використовувати паралельні обчислення для зменшення часу навчання.
- Стійкість до шуму: Бутстреп-агрегація може бути стійким до шуму в даних, оскільки вона використовує середнє або голосування з кількох моделей. Вона може зменшити вплив випадкових помилок або викидів на кінцевий прогноз.

Ураховуючи ці переваги та недоліки, бутстреп-агрегація є потужним методом ансамблевого навчання, який може поліпшити якість прогнозів і допомогти уникнути перенавчання. Вона широко використовується в різних областях, включаючи машинне навчання, статистику і прогнозування.

2. Бустінг (Boosting): цей метод базується на послідовному навчанні слабких моделей навчання (наприклад, дерева рішень), при цьому кожна наступна модель намагається скоригувати помилки попередньої моделі.

Переваги методу бустінгу:

- Висока точність: Бустінг зазвичай дає дуже точні результати, оскільки кожна наступна модель намагається скоригувати помилки попередньої моделі. Це дозволяє досягти високої точності прогнозів.
- Здатність до моделювання складних залежностей: Бустінг може моделювати складні залежності між змінними, оскільки він може використовувати композицію багатьох слабких моделей, які разом утворюють потужну прогностичну модель.

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		18

обчислювальних ресурсів та часу, особливо якщо використовуються складні моделі або великі набори даних.

Загалом, стекінг є потужним засобом ансамблевого навчання, який може покращити точність прогнозування і моделювати складні залежності в даних. Однак, використання стекінгу вимагає налаштування моделей та може призвести до підвищеного ризику перенавчання. Також, використання стекінгу може бути обчислювально витратним процесом. Якщо належним чином налаштувати та застосувати стекінг, він може бути потужним інструментом для покращення прогнозування і отримання кращих результатів.

У загальному, ансамбль методів машинного навчання може допомогти підвищити точність прогнозування в порівнянні з окремими моделями. Бутстреп-агрегація (Bagging) може знизити варіацію та покращити стійкість моделі до перенавчання, тоді як бустінг (Boosting) може покращити точність шляхом зменшення помилок та підсилення слабких моделей навчання. Однак, бустінг може бути чутливим до шуму в даних та перенавчання, тому необхідно враховувати ці фактори при виборі підходів для ансамблю. Стекінг (Stacking) може покращити точність, проте він вимагає більш складного процесу навчання та може стати більш обчислювально витратним. Враховуючи ці фактори, бустінг може бути одним з найбільш ефективних підходів до ансамблю методів машинного навчання.

Існує кілька алгоритмів бустингу, які використовуються для підвищення точності моделей машинного навчання. Основні алгоритми бустингу включають:

1. AdaBoost (Adaptive Boosting): Цей алгоритм навчає послідовні слабкі класифікатори на даних, з фокусом на екземплярах, які були неправильно класифіковані попередніми моделями. Важливість кожного екземпляра змінюється залежно від того, наскільки часто він був неправильно класифікований, дозволяючи наступним моделям навчатися на важчих екземплярах.

										ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата							21

6. Розробка алгоритму валідації моделей для перевірки точності прогнозування

7. Реалізація та валідація ансамблевого навчання моделей для прогнозування врожайності ягід

Завдання 1-3 розглянуто в параграфах 1.1 - 1.4 відповідно, а виконання завдань 3-6 представлені у параграфах 2.1-3.3.

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		25

XGBoost також обробляє відсутні дані, зменшуючи необхідність передобробки даних.

XGBoost (eXtreme Gradient Boosting) є високоефективною реалізацією алгоритму градієнтного бустінгу. На відміну від традиційного градієнтного бустінгу, XGBoost включає регуляризацію в свою функцію втрат, що допомагає контролювати перенавчання.

Для пояснення математичного аспекту, давайте розглянемо функцію втрат XGBoost:

$$L(\varphi) = \sum l(y^i, y^i) + \sum \Omega(fj),$$

де:

$l(y^i, y^i)$ - це функція втрат, що вимірює різницю між спостережуваною міткою y^i та прогнозованою міткою y^i .

$\Omega(fj)$ - це регуляризація, що штрафує складність моделі. Це може бути, наприклад, L1 або L2 регуляризація.

Функція втрат зазвичай є квадратичною втратою для задач регресії і логістичною втратою для задач класифікації. Однак XGBoost дозволяє використовувати користувацькі функції втрат.

XGBoost намагається мінімізувати цю функцію втрат шляхом додавання нових моделей, які найкращим чином зменшують втрату (тобто виправляють помилки попередніх моделей), одночасно контролюючи складність моделі через регуляризацію.

LightGBM (Light Gradient Boosting Machine) - це алгоритм, розроблений Microsoft, який використовує техніку називається градієнтний бустінг з основою на деревах (GOSS). LightGBM здатний обробляти великі набори даних та має високу швидкість обчислень.

LightGBM (Light Gradient Boosting Machine) є алгоритмом градієнтного бустінгу, який використовує дерево, основане на градієнтах, для ефективного навчання. Він використовує дві ключові техніки для зменшення

									Арк.
									27
Зм.	Арк.	№ докум.	Підпис	Дата					

ДП.КН. 8351441.062ПЗ

за допомогою Optuna. Потім створюються моделі з найкращими гіперпараметрами, вибраними Optuna. Далі проводиться відбір ознак за допомогою методу SelectKBest. На основі вибраних ознак тренується ансамбль моделей. Нарешті, цей ансамбль використовується для прогнозування на тестовому наборі даних, і його продуктивність оцінюється за допомогою коефіцієнта детермінації R^2 .

Алгоритм навчання моделей представлено кроками та блок-схемою (див.дод.А):

Крок 1. Попередня обробка даних: Дані були розбиті на незалежні (X) та залежні (y) змінні. Змінна 'yield' була виділена як цільова змінна. Далі дані були розбиті на набори для тренування та тестування в співвідношенні 80:20. Нарешті, дані були стандартизовані за допомогою StandardScaler.

Крок 2. Визначення цілей Optuna для кожної моделі: Для кожної моделі (XGBoost, LightGBM, CatBoost) було визначено цільову функцію для оптимізації гіперпараметрів. Ці функції включали створення моделі, її тренування на навчальних даних та оцінку за допомогою коефіцієнта детермінації R^2 .

Крок 3. Оптимізація гіперпараметрів з Optuna: Для кожної моделі було запущено процес оптимізації гіперпараметрів Optuna. Optuna використовувала випадковий пошук для вибору гіперпараметрів, які максимізують R^2 .

Крок 4. Виведення найкращих гіперпараметрів: Для кожної моделі було виведено найкращі гіперпараметри, знайдені в процесі оптимізації Optuna.

Крок 5. Тренування базових моделей з найкращими параметрами: За допомогою найкращих гіперпараметрів було створено та натреновано три моделі: XGBoost, LightGBM та CatBoost.

Крок 6. Вибір ознак та створення ансамблю моделей: Було використано SelectKBest для вибору 10 найкращих ознак. Потім було

										ДП.КН. 8351441.062ПЗ	Арк.
											30
Зм.	Арк.	№ докум.	Підпис	Дата							

створено ансамбль моделей за допомогою VotingRegressor, що об'єднує прогнози від трьох базових моделей.

Крок 7. Тренування ансамблю моделей: Нарешті, ансамбль моделей було натреновано на вибраних ознаках тренувального набору даних.

Алгоритм навчання моделей прогнозування, може бути особливо корисним для створення інтелектуального модулю прогнозування врожайності ягід. Використання ансамблю трьох різних моделей машинного навчання дозволяє використовувати переваги кожної з них, що може призвести до більш точних прогнозів. Оптимізація гіперпараметрів з Optuna допомагає забезпечити, що кожна модель налаштована таким чином, щоб найкращим чином використовувати доступні дані. Відбір ознак за допомогою SelectKBest може допомогти зосередитися на найбільш важливих ознаках, що допомагає покращити продуктивність моделі, зменшуючи при цьому обсяг обчислень.

2.3 Алгоритм валідації моделей прогнозування врожайності ягід

Алгоритм розроблений для проведення валідації та тестування навченої ансамблевої моделі для прогнозування врожайності ягід. Він має на меті перевірити ефективність моделі за допомогою різних методів, включаючи перехресну валідацію, обчислення середньоквадратичного відхилення (MSE) та коефіцієнта детермінації (R-квадрат), а також аналіз залишків.

Також включені кроки, які забезпечують коректну передобробку нових даних, прогнозування за допомогою моделі та подальше збереження та візуалізацію результатів. Це не тільки дає змогу оцінити якість прогнозування, але й надає інформацію, яка може бути використана для подальшого вдосконалення моделі.

Цей алгоритм має широкий спектр застосувань, особливо в галузях, де важливо точне прогнозування, таких як агрономія, екологія та інші сфери прикладних наук.

Отже, алгоритм описаний послідовністю кроків та блок-схемою (див.дод.Б):

Вхід - Навчений ансамбль моделей на вибраних характеристиках.

Крок 1. Проведіть перехресну валідацію ансамблю моделей, використовуючи KFold, та обчисліть R-квадрат для оцінки ефективності моделі.

KFold, або K-кратне перехресне перевірку, це метод для оцінки моделі машинного навчання. Основна ідея полягає в тому, що навчальний набір даних ділиться на K рівних піднаборів. Потім проводиться K ітерацій навчання та тестування, кожен раз використовуючи один з піднаборів як тестовий набір, а решту - як навчальний.

Для кожної ітерації i (де i знаходиться в діапазоні від 1 до K), обчислюється помилка моделі на тестовому піднаборі i . Потім ці помилки усереднюються, щоб отримати загальну оцінку моделі.

Математично це можна описати наступним чином:

Розбиття набору даних на K піднаборів:

$$S_1, S_2, \dots, S_k$$

Для кожного піднабору S_i :

а. Використовуйте S_i як тестовий набір, а всі інші піднабори як навчальний набір.

б. Навчіть модель на навчальному наборі та обчисліть помилку на тестовому наборі.

Усередніть помилки з усіх K ітерацій, щоб отримати загальну оцінку моделі.

Це дає уявлення про те, як модель буде працювати на нових даних, що є важливим для визначення її загальної ефективності.

									Арк.
									32
Зм.	Арк.	№ докум.	Підпис	Дата					

R-квадрат, або коефіцієнт детермінації, це статистична міра, яка використовується для оцінки якості моделі в регресійному аналізі. Він показує, яку частку загальної дисперсії залежної змінної можна пояснити за допомогою моделі.

R-квадрат вимірюється від 0 до 1, де 1 означає, що модель повністю пояснює дисперсію залежної змінної, а 0 означає, що модель не пояснює дисперсію залежної змінної взагалі.

Математично R-квадрат можна виразити наступним чином:

$$R^2 = 1 - \left(\frac{SS_{res}}{SS_{tot}} \right),$$

де

- SS_{res} є сумою квадратів залишків (різниць між фактичними та прогнозованими значеннями), і
- SS_{tot} є сумою квадратів відносно середнього значення залежної змінної.

Залишкова сума квадратів (SS_{res}) вимірює внутрішню варіативність, яку модель не змогла виявити, а загальна сума квадратів (SS_{tot}) вимірює загальну варіативність. Тому R-квадрат відображає частку загальної варіативності, яку модель змогла виявити.

Крок 2. Протестуйте ансамбль моделей на тестовій вибірці та обчисліть середньоквадратичну помилку (MSE) та R-квадрат для оцінки ефективності моделі.

Середньоквадратична помилка (MSE - Mean Squared Error) - це часто використовувана міра відмінності між значеннями прогнозованими моделлю та справжніми значеннями. Вона вимірює середній квадрат відхилень між цими двома наборами значень.

Математично MSE можна виразити таким чином:

$$MSE = \left(\frac{1}{n} \right) * \Sigma(actual - prediction)^2$$

									ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата						33

де:

- n є кількістю точок даних
- $actual$ є справжнім значенням
- $prediction$ є прогнозованим значенням
- Σ означає суму по всіх точках даних.

Основна ідея полягає в тому, щоб підрахувати квадрат відхилення для кожної точки даних, а потім обчислити середнє цих квадратів. Більш високі значення MSE відповідають більшій помилці прогнозування, тобто модель з більшим значенням MSE виконує прогноз гірше, ніж модель з меншим значенням MSE.

Крок 3. Побудуйте графік залишків для визначення будь-яких систематичних відхилень у прогнозах моделі.

Крок 4. Застосуйте модель до нових даних, використовуючи ті ж самі кроки передобробки даних, що були використані для навчальної вибірки.

Крок 5. Зробіть прогнози для нових даних, використовуючи вашу модель.

Крок 6. Створіть новий DataFrame з прогнозованими значеннями і збережіть його у файл csv.

Крок 7. Візуалізуйте прогнозовані результати за допомогою графіків або діаграм для кращого розуміння результатів.

Алгоритм валідації та тестування, представлений в цій роботі, виявився ефективним у демонстрації роботи інтелектуального модуля прогнозування врожайності ягід. Він дозволив виміряти точність прогнозів моделі та її здатність узагальнювати на нові дані.

Перехресна валідація та обчислення таких метрик як середньоквадратичне відхилення (MSE) та коефіцієнт детермінації (R-квадрат) допомогли оцінити ефективність ансамблевої моделі. Аналіз залишків виявив важливі відомості про прогнози моделі та її здатність адекватно моделювати відношення між характеристиками та цільовою змінною.

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		34

**3 ПРОГРАМНО-ТЕХНОЛОГІЧНЕ ЗАБЕЗПЕЧЕННЯ
ІНТЕЛЕКТУАЛЬНОГО МОДУЛЯ ПРОГНОЗУВАННЯ ВРОЖАЙНОСТІ
ЯГІД**

3.1 Розвідувальний аналіз набору даних

Набір даних, який будемо використовувати для прогностичного моделювання, був створений за допомогою моделі симуляції опилення дикої чорниці Wild Blueberry Pollination Simulation Model [47]. Це відкрита, просторово-експліцитна комп'ютерна програма, яка дозволяє досліджувати, як різні фактори, включаючи просторове розташування рослин, кроспилення та самоопилення, склад видів бджіл та погодні умови, поодиночі та в комбінації, впливають на ефективність опилення та врожай агроecosистеми дикої чорниці. Модель симуляції була підтверджена полевыми спостереженнями та експериментальними даними, зібраними в штаті Мен, США, та в Канадських Марітимах протягом останніх 30 років, і тепер є корисним інструментом для перевірки гіпотез та розвитку теорії досліджень опилення дикої чорниці.

Отже, проведемо розвідувальний аналіз цього набору даних. Цей набір даних включає 15 289 записів (Таблиця 3.1). Для кожної з характеристик зазначено середнє значення (mean), стандартне відхилення (std), мінімальне (min), максимальне (max) значення, а також 25%, 50% і 75% перцентилі (відповідно 1-й кuartиль, медіана і 3-й кuartиль).

1. Розмір клона варіюється від 10 до 40 м², з середнім значенням 19.7 м².
2. Плотність бджіл в полі має різний розподіл для різних видів бджіл. Наприклад, середнє значення для медоносних бджіл становить 0.39 бджол/м²/хв, для шмелів - 0.29, для *Andrena* - 0.49, для *Osmia* - 0.59.
3. Діапазони температур варіюють від 69.7 до 94.6 °C для верхнього діапазону і від 50.2 до 68.2 °C для нижнього діапазону. Середні значення становлять відповідно 68.7 °C та 48.6 °C.

						<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата			35

4. Кількість днів з дощем варіюється від 1 до 34 днів, із середнім значенням 18.7 днів.
5. Середнє значення днів з дощем протягом всього періоду цвітіння становить 0.32 дня.
6. Значення показників "fruitset", "fruitmass", "seeds" та "yield" варіюються відповідно від 0.19 до 0.65, від 0.31 до 0.54, від 22.08 до 46.59 та від 1945.53 до 8969.4 з середніми значеннями 0.50, 0.45, 36.16 та 6025.19.

Ця інформація важлива для розуміння розподілу даних і може допомогти у виявленні аномалій та визначенні стратегії обробки даних.

Далі визначимо кореляцію (Рис.3.1) між різними характеристиками і урожайністю. Значення кореляції можуть варіюватися від -1 до 1, де -1 позначає повну негативну кореляцію, 1 - повну позитивну кореляцію, а 0 означає відсутність кореляції.

Розмір клона ("clonesize"), кількість медоносних бджіл ("honeybee") та всі температурні показники ("MaxOfUpperTRange", "MinOfUpperTRange", "AverageOfUpperTRange", "MaxOfLowerTRange", "MinOfLowerTRange", "AverageOfLowerTRange") мають негативну кореляцію з урожайністю. Це означає, що зі збільшенням цих показників урожайність зменшується. Найсильнішу негативну кореляцію мають "AverageRainingDays" (-0.48) і "RainingDays" (-0.47).

Шмелі ("bumbles"), *Andrena* ("andrena") та *Osmia* ("osmia") мають позитивну кореляцію з урожайністю, що означає, що зі збільшенням цих показників урожайність також збільшується.

Таблиця 3.1. Описова статистика набору даних

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		36

	clonesize	honeybee	bumbles	andrena	osmia
count	15289.0	15289.0	15289.0	15289.0	15289.0
mean	19.7	0.4	0.3	0.5	0.6
std	6.6	0.4	0.1	0.1	0.1
min	10.0	0.0	0.0	0.0	0.0
25%	12.5	0.2	0.2	0.4	0.5
50%	25.0	0.5	0.2	0.5	0.6
75%	25.0	0.5	0.4	0.6	0.8
max	40.0	18.4	0.6	0.8	0.8

	MaxOfUpper TRange	MinOfUpper TRange	AverageOfUpper TRange	MaxOfLower TRange	MinOfLower TRange	AverageOfLower TRange
--	----------------------	----------------------	--------------------------	----------------------	----------------------	--------------------------

	15289.0	15289.0	15289.0	15289.0	15289.0	15289.0
--	---------	---------	---------	---------	---------	---------

count	82.2	49.7	68.7	59.2	28.7	48.6
mean	9.1	5.5	7.6	6.6	3.2	5.4
std	69.7	39.0	58.2	50.2	24.3	41.2
min	77.4	46.8	64.7	55.8	27.0	45.8
25%	86.0	52.0	71.9	62.0	30.0	50.8
50%	86.0	52.0	71.9	62.0	30.0	50.8
75%	94.6	57.2	79.0	68.2	33.0	55.9
max						

	Raining Days	AverageRaining Days	fruitset	fruitmass	seeds	yield
--	-----------------	------------------------	----------	-----------	-------	-------

count	15289.0	15289.0	15289.0	15289.0	15289.0	15289.0
mean	18.7	0.3	0.5	0.4	36.2	6025.2
std	11.7	0.2	0.1	0.0	4.0	1337.1
min	1.0	0.1	0.2	0.3	22.1	1945.5
25%	16.0	0.3	0.5	0.4	33.2	5128.2
50%	16.0	0.3	0.5	0.4	36.0	6117.5
75%	24.0	0.4	0.6	0.5	39.2	7019.7
max	34.0	0.6	0.7	0.5	46.6	8969.4

Найсильнішу позитивну кореляцію з урожайністю мають "fruitset" (0.89), "seeds" (0.87) та "fruitmass" (0.83). Це означає, що ці показники найбільше впливають на урожайність.

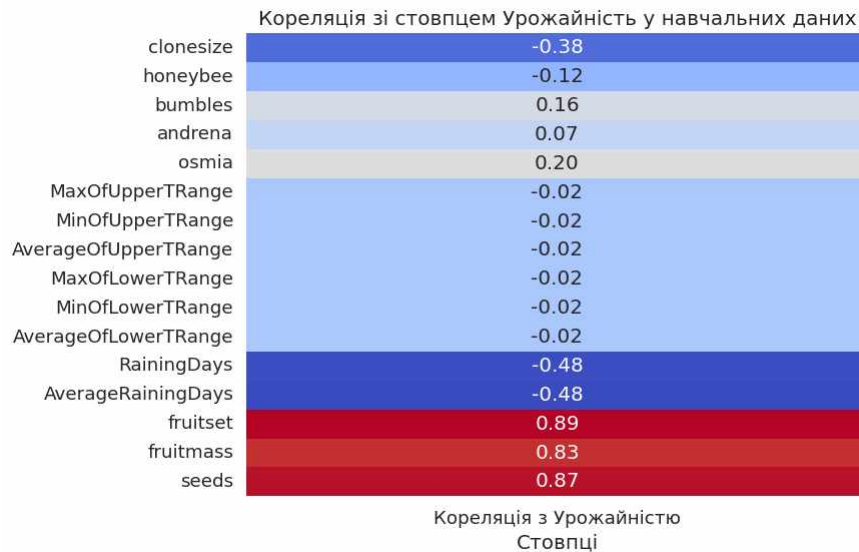


Рисунок 3.1 - Теплова матриця кореляції

Отже, виведемо:

Важливі характеристики на основі порогу кореляції зі стовпчиком врожайності плодів: ['clonesize', 'honeybee', 'bumbles', 'osmia', 'RainingDays', 'AverageRainingDays', 'fruitset', 'fruitmass', 'seeds', 'yield']

Далі проведемо аналіз результатів діаграм віолончелі (Рис. 3.2), які дають уявлення про розподіл залежності між важливими ознаками і цільовою змінною "урожайність". Діаграма віолончелі показує густину розподілу змінної "урожайність" для кожного значення важливої ознаки.

За допомогою цих діаграм можна проаналізувати, як залежність між важливими ознаками та урожайністю може варіюватись. Якщо графік діаграми віолончелі вузький у верхній частині та широкий у нижній частині, це означає, що вищі значення ознаки мають тенденцію бути пов'язаними з вищими значеннями урожайності. Зворотна ситуація спостерігається, якщо графік діаграми вузький у нижній частині та широкий у верхній частині.

Отже, проведено розвідувальний аналіз набору даних, який було згенеровано за допомогою моделі симуляції опилення дикої чорниці (Wild Blueberry Pollination Simulation Model). Модель була перевірена та валідована за допомогою польових спостережень та експериментальних даних, зібраних у штаті Мен (США) та Канадських Марітимах протягом останніх 30 років.

Розглянули описові статистики для кожної з ознак набору даних, які включали середнє значення, стандартне відхилення, мінімальне та максимальне значення, а також кватильні значення.

Далі проаналізували кореляцію між різними характеристиками та урожайністю. Виявили, що певні ознаки, такі як "fruitset", "seeds" та "fruitmass", мають сильну позитивну кореляцію з урожайністю, в той час як інші, як-от "AverageRainingDays" та "RainingDays", мають сильну негативну кореляцію.

Ці аналізи є важливими для розуміння взаємозв'язків між ознаками та цільовим показником - у цьому випадку, урожайністю. Ці знання можуть бути корисними для створення та покращення прогностичних моделей, які передбачають урожайність на основі вхідних характеристик.

3.2 Ансамблеве навчання моделей

У цьому параграфі проведемо ансамблеве навчання моделей. Для цього спершу дані будуть очищені, розділені на навчальні та тестові набори, а також нормалізовані для оптимальної роботи моделі машинного навчання. Це важливі кроки, які допомагають забезпечити найкращу можливу продуктивність нашої моделі.

```
1 # Preprocess the data
2 X = data.drop('yield', axis=1)
3 y = data['yield']

1 # Split the data into training and testing sets
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
3
4 # Scale the data
5 scaler = StandardScaler()
6 X_train = scaler.fit_transform(X_train)
7 X_test = scaler.transform(X_test)
```

									Арк.
									40
Зм.	Арк.	№ докум.	Підпис	Дата					

ДП.КН. 8351441.062ПЗ

Цей код виконує кілька ключових кроків для підготовки даних для машинного навчання:

1. Попередня обробка даних: В цьому випадку, ми розділяємо наш набір даних на цільову змінну 'yield' та інші характеристики (або ознаки). Це робиться за допомогою методу drop в pandas DataFrame, який видаляє стовпець 'yield' з основного набору даних data і зберігає інші стовпці в X. Значення 'yield' зберігаються в y.
2. Розбиття даних на навчальний та тестовий набори: Далі ми розбиваємо наші дані на два різних набори - навчальний і тестовий. Це допомагає нам перевірити, наскільки добре наша модель працює, коли вона отримує нові, раніше невідомі їй дані. Ми використовуємо 80% даних для навчання і 20% для тестування.
3. Масштабування даних: В кінці ми масштабуємо наші дані, використовуючи StandardScaler від sklearn. Це перетворює всі наші числові ознаки так, що вони матимуть середнє значення 0 і стандартне відхилення 1. Це дуже важливий крок, особливо для моделей, які вимагають, щоб усі ознаки мали подібний масштаб.

```
def xgb_objective(trial):
    n_estimators = trial.suggest_int("n_estimators", 100, 500)
    learning_rate = trial.suggest_float("learning_rate", 0.01, 0.2)
    max_depth = trial.suggest_int("max_depth", 3, 10)
    subsample = trial.suggest_float("subsample", 0.5, 1)
    colsample_bytree = trial.suggest_float("colsample_bytree", 0.5, 1)

    model = XGBRegressor(
        n_estimators=n_estimators,
        learning_rate=learning_rate,
        max_depth=max_depth,
        subsample=subsample,
        colsample_bytree=colsample_bytree,
        random_state=42
    )

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)
    r2 = r2_score(y_test, y_pred)
    return r2
```

									ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата						41

```

def lgbm_objective(trial):
    n_estimators = trial.suggest_int("n_estimators", 100, 500)
    learning_rate = trial.suggest_float("learning_rate", 0.01, 0.2)
    max_depth = trial.suggest_int("max_depth", 3, 10)
    num_leaves = trial.suggest_int("num_leaves", 31, 127)
    subsample = trial.suggest_float("subsample", 0.5, 1)
    colsample_bytree = trial.suggest_float("colsample_bytree", 0.5, 1)

    model = LGBMRegressor(
        n_estimators=n_estimators,
        learning_rate=learning_rate,
        max_depth=max_depth,
        num_leaves=num_leaves,
        subsample=subsample,
        colsample_bytree=colsample_bytree,
        random_state=42
    )

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)
    r2 = r2_score(y_test, y_pred)
    return r2

```

```

def catboost_objective(trial):
    iterations = trial.suggest_int("iterations", 100, 500)
    learning_rate = trial.suggest_float("learning_rate", 0.01, 0.2)
    depth = trial.suggest_int("depth", 3, 10)
    l2_leaf_reg = trial.suggest_int("l2_leaf_reg", 1, 9)
    subsample = trial.suggest_float("subsample", 0.5, 1)

    model = CatBoostRegressor(
        iterations=iterations,
        learning_rate=learning_rate,
        depth=depth,
        l2_leaf_reg=l2_leaf_reg,
        subsample=subsample,
        random_state=42,
        verbose=0
    )

    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)
    r2 = r2_score(y_test, y_pred)
    return r2

```

Цей код визначає три функції об'єктиву для оптимізації гіперпараметрів трьох моделей машинного навчання: XGBoost, LightGBM і CatBoost за допомогою Optuna. Optuna - це бібліотека Python, що використовується для автоматичної оптимізації гіперпараметрів.

1. **xgb_objective**: Ця функція визначає об'єктив для оптимізації моделі XGBoost. Вона використовує методи `suggest_int` і `suggest_float` Optuna

											ДП.КН. 8351441.062ПЗ	Арк.
												42
Зм.	Арк.	№ докум.	Підпис	Дата								

Цей код використовує Optuna для оптимізації гіперпараметрів трьох моделей машинного навчання: XGBoost, LightGBM та CatBoost. Це робиться за допомогою оптимізаційних функцій, визначених в попередньому коді.

1. Створення вибірника: Використовується вибірник TPE (Tree-structured Parzen Estimator), який є методом оптимізації гіперпараметрів, що базується на байєсівському підході. Цей вибірник забезпечує рандомізацію, але вихідні дані він генерує на основі попередніх результатів, щоб знайти найкращий набір гіперпараметрів.
2. Створення та оптимізація досліджень: Для кожної з трьох моделей створюється окреме дослідження за допомогою `optuna.create_study`. Напрямок "maximize" вказує, що метою є максимізація результату функції об'єктиву. Далі ці дослідження оптимізуються за допомогою відповідних функцій об'єктиву, визначених раніше. Кількість випробувань встановлюється як 50 для кожної моделі.
3. Виведення найкращих гіперпараметрів: Після завершення оптимізації для кожної моделі виводяться найкращі знайдені гіперпараметри. Це дозволяє побачити, які значення гіперпараметрів привели до найкращого результату для кожної моделі.

В результаті отримуємо найкращі гіперпараметри для кожної з трьох моделей, що були оптимізовані:

```
Найкращі параметри для XGBoost від Optuna: {'n_estimators': 291, 'learning_rate': 0.02572640965568375, 'max_depth': 4, 'subsample': 0.552187345043806, 'colsample_bytree': 0.8272028535549523}
```

```
Найкращі параметри для LightGBM від Optuna: {'n_estimators': 217, 'learning_rate': 0.07053452738526493, 'max_depth': 4, 'num_leaves': 86, 'subsample': 0.9281322440018367, 'colsample_bytree': 0.6325640166339597}
```

```
Найкращі параметри для CatBoost від Optuna: {'iterations': 288, 'learning_rate': 0.0537104881625418, 'depth': 9, 'l2_leaf_reg': 2, 'subsample': 0.542177159163148}
```

XGBoost: Optuna знайшла найкращі гіперпараметри для моделі XGBoost. Кількість дерев (`n_estimators`) є 291, швидкість навчання (`learning_rate`) становить приблизно 0.026, максимальна глибина дерева

										Арк.
										44
Зм.	Арк.	№ докум.	Підпис	Дата						

ДП.КН. 8351441.062ПЗ

Отже, було реалізовано дві ключові операції: вибір найважливіших ознак за допомогою SelectKBest, а також створення ансамблю моделей на основі XGBoost, LightGBM та CatBoost. Вибір ознак дозволяє зосередитися на найінформативніших змінних, що сприяє покращенню якості прогнозування. Ансамбль моделей об'єднує сили трьох потужних алгоритмів машинного навчання, що призводить до збільшення точності прогнозів. Усі ці операції разом створюють сильну основу для інтелектуального модулю прогнозування врожайності ягід.

3.3 Валідація моделей

Далі проведемо перехресну перевірку (k-fold cross-validation) для оцінки ефективності ансамблевої моделі прогнозування.

```
1 kfold = KFold(n_splits=5, shuffle=True, random_state=42)
2 cv_scores = cross_val_score(ensemble, X_train_selected, y_train, cv=kfold, scoring='r2')
3 print("Оцінки перехресної перевірки: ", cv_scores)
4 print("Середнє R-квадрат перехресної перевірки: {:.2f}".format(np.mean(cv_scores)))
5 print("Стандартне відхилення R-квадрат перехресної перевірки: {:.2f}".format(np.std(cv_scores)))
```

Оцінки перехресної перевірки: [0.79021216 0.82459323 0.82172327 0.83670101 0.81620097]
Середнє R-квадрат перехресної перевірки: 0.82
Стандартне відхилення R-квадрат перехресної перевірки: 0.02

У цьому випадку використовується 5-кратна перехресна перевірка, що означає, що набір тренувальних даних поділяється на 5 частин, і модель навчається та тестується 5 разів, кожен раз з використанням різного піднабору даних для тестування.

Оцінки R-квадрат, отримані в результаті перехресної перевірки, вказують на досить високу здатність моделі прогнозувати врожайність ягід. Середнє значення R-квадрат становить 0.82, що означає, що модель може пояснити приблизно 82% варіативності врожайності ягід. Це свідчить про досить високий рівень точності прогнозування.

Стандартне відхилення R-квадрат, що дорівнює 0.02, вказує на те, що результати перехресної перевірки досить стабільні. Модель показує схожі

										ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата							48

результати на різних підмножинах тренувальних даних, що свідчить про її надійність і стабільність.

У цілому, отримані результати свідчать про те, що ансамблева модель з гіперпараметрами, отриманими в результаті оптимізації Optuna, є досить ефективною для прогнозування врожайності ягід.

```
2 y_pred = ensemble.predict(X_test_selected)
3
4 mse = mean_squared_error(y_test, y_pred)
5 r2 = r2_score(y_test, y_pred)
6
7 print("Середнє квадратичне відхилення: {:.2f}".format(mse))
8 print("Середнє квадратичне відхилення: {:.2f}".format(r2))
```

Середнє квадратичне відхилення: 311929.68
Середнє квадратичне відхилення: 0.82

Цей код використовує ансамблеву модель для прогнозування врожайності ягід на основі тестового набору даних. Він також вимірює дві ключові метрики оцінки ефективності моделі: середньоквадратичне відхилення (MSE) та коефіцієнт детермінації R-квадрат.

Середнє квадратичне відхилення (MSE): Це метрика, яка вимірює середню квадратичну різницю між фактичними та прогнозованими значеннями. Воно є одним із найбільш поширених критеріїв оцінки для завдань регресії. В даному випадку, MSE становить 311929.68, що є відносно великим значенням. Це може свідчити про те, що модель допускає помітні помилки прогнозування, але конкретна інтерпретація залежить від масштабу врожайності ягід та її варіативності в даних.

Коефіцієнт детермінації (R-квадрат): Ця метрика вимірює долю варіативності відгуку, яку модель може пояснити. Його значення варіюється від 0 до 1, де 1 означає, що модель може пояснити всю варіативність відгуку, а 0 - що модель не може пояснити ніякої варіативності. В даному випадку, R-квадрат становить 0.82, що свідчить про те, що модель може пояснити 82% варіативності врожайності ягід. Це є досить високим значенням, що свідчить про ефективність моделі.

									ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата						49

Графік (Рис.3.5) залишкових помилок важливий для оцінки якості моделі. Залишки розподілені випадковим чином навколо нуля, що свідчило б про відсутність систематичних помилок.

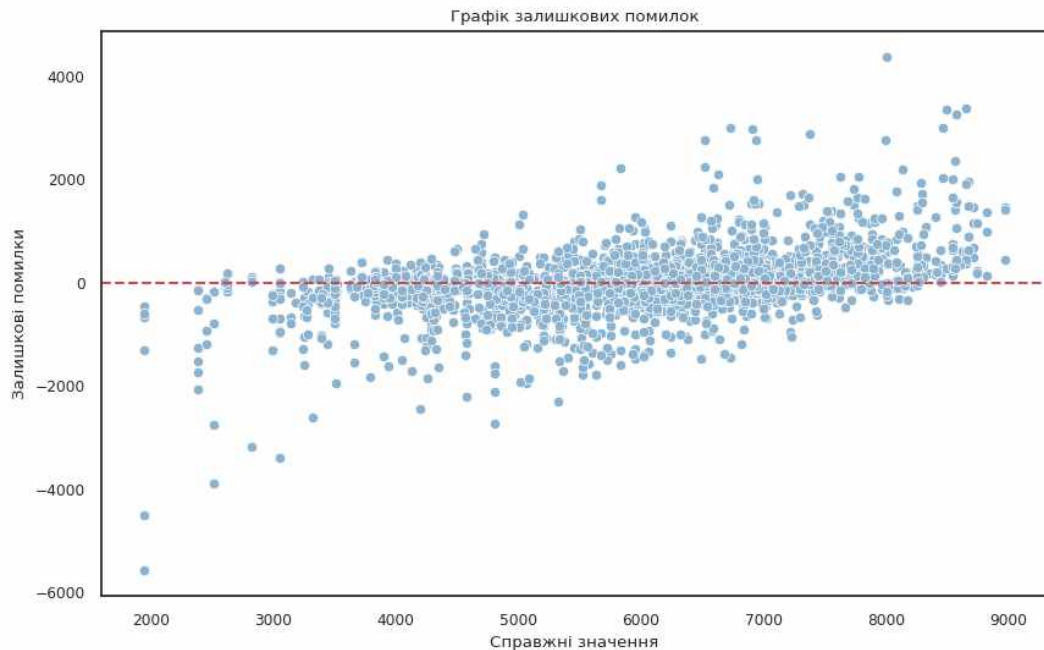


Рисунок 3.5 - Графік залишкових помилок

Далі здійснюємо інженерію ознак (feature engineering) на тестовому наборі даних (test_df), створюючи нові ознаки на основі наявних.

```

1 # 1. Temperature range
2 test_df['TemperatureRange'] = test_df['MaxOfUpperTRange'] - test_df['MinOfLowerTRange']
3
4 # 2. Temperature extremes
5 threshold_high = 72
6 threshold_low = 50
7
8 test_df['ExtremeHighTemp'] = (test_df['AverageOfUpperTRange'] > threshold_high).astype(int)
9 test_df['ExtremeLowTemp'] = (test_df['AverageOfLowerTRange'] < threshold_low).astype(int)
10
11 # 3. Total bee density
12 test_df['TotalBeeDensity'] = test_df['honeybee'] + test_df['bumbles'] + test_df['andrena'] + test_df['osmia']
13
14 # 4. Bee species dominance
15 total_density = test_df['honeybee'] + test_df['bumbles'] + test_df['andrena'] + test_df['osmia']
16 test_df['HoneybeeDominance'] = test_df['honeybee'] / total_density
17 test_df['BumblesBeeDominance'] = test_df['bumbles'] / total_density
18 test_df['AndrenaBeeDominance'] = test_df['andrena'] / total_density
19 test_df['OsmiaBeeDominance'] = test_df['osmia'] / total_density
20
21 # 5. Rain intensity
22 test_df['RainIntensity'] = test_df['AverageRainingDays'] / test_df['RainingDays']
23
24 # 6. Interaction features
25 test_df['BeeDensity_TemperatureInteraction'] = test_df['TotalBeeDensity'] * test_df['TemperatureRange']
26 test_df['BeeDensity_RainInteraction'] = test_df['TotalBeeDensity'] * test_df['RainIntensity']

```

TemperatureRange: Ця нова ознака представляє діапазон температури, який обчислюється як різниця між максимальною та мінімальною температурою.

З результатів видно (Рис.3.7), що тренд прогнозу відповідає реальним значенням, це означає, що запропоновані моделі доволі не погано прогнозують результати.

← predictions.csv		
	A	B
1	id	yield
2	15289	4317.541124
3	15290	6236.049772
4	15291	7242.849531
5	15292	4842.237378
6	15293	3608.927953
7	15294	5113.702693
8	15295	7282.286407
9	15296	6529.989453
10	15297	8095.250667
11	15298	4303.136139
12	15299	5805.145707
13	15300	5791.312897
14	15301	5565.612833
15	15302	5674.574637
16	15303	4956.173162
17	15304	4666.273687
18	15305	4677.971637

Рисунок 3.6 - Вивід результатів 'predictions.csv'

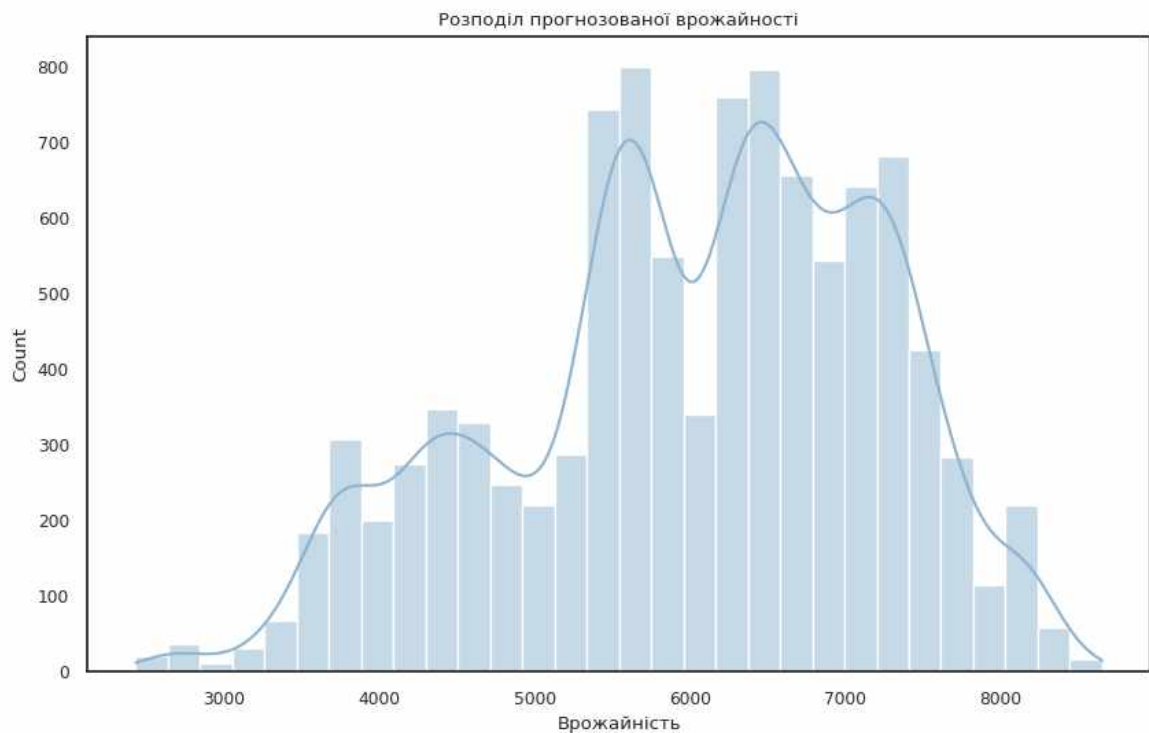


Рисунок 3.7 - Вивід результатів у вигляді графіка

Зм.	Арк.	№ докум.	Підпис	Дата

ДП.КН. 8351441.062ПЗ

Арк.

53

Інтелектуальний модуль прогнозування, що використовує ансамбль моделей XGBoost, LightGBM та CatBoost, успішно впорався з завданням. Ці моделі були натреновані на основі відповідних гіперпараметрів, знайдених за допомогою оптимізації Optuna, що дозволило збільшити їхню продуктивність.

Спостерігаючи за графіками прогнозів та реальних значень, можна бачити, що вони мають подібні тенденції, що вказує на високу здатність моделей передбачити врожайність.

Це підтверджує, що використання машинного навчання та оптимізації гіперпараметрів може бути дуже ефективним для прогнозування врожайності ягід. Модель може бути використана фермерами або агрономами для оптимізації вирощування ягід та планування врожаю.

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		54

ВИСНОВКИ

Використання бустингових алгоритмів, таких як XGBoost, LightGBM та CatBoost, у ролі основи для інтелектуального модулю прогнозування врожайності ягід, демонструє значний потенціал. Ці алгоритми надають розробникам потужні інструменти для створення точних прогнозів.

Використання цих алгоритмів в якості ансамблю моделей дозволяє використовувати сильні сторони кожного з них, щоб компенсувати можливі слабкі місця. Це забезпечує найкращу можливу точність прогнозування, що стає важливим інструментом для прийняття обґрунтованих рішень в агробізнесі.

Розроблений алгоритм навчання моделей прогнозування врожайності ягід представляє собою важливий інструмент для створення інтелектуального модуля прогнозування. Використання комбінації трьох різних моделей машинного навчання забезпечує можливість максимально ефективного використання їх сильних сторін для отримання більш точних прогнозів. Налаштування гіперпараметрів за допомогою Optuna гарантує, що кожна модель оптимально адаптована до конкретного набору даних. Застосування методу SelectKBest для відбору ознак сприяє концентрації на найважливіших з них, що підвищує ефективність моделі та зменшує обсяг необхідних обчислень.

Розроблений алгоритм валідації та тестування виявив свою ефективність у визначенні роботоздатності інтелектуального модулю прогнозування врожайності ягід. Цей процес надав можливість оцінити точність моделі та її спроможність генералізувати відомості на нові дані. Імплементация крос-валідації та використання метрик як-от середньоквадратичне відхилення (MSE) та коефіцієнт детермінації (R-квадрат) сприяли вимірюванню продуктивності ансамблевої моделі. Додатково, аналіз залишків надав цінні відомості про прогнозні здібності моделі та її здатність адекватно відобразити відносини між ознаками та цільовою змінною.

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		55

В ході роботи було проведено аналіз даних, згенерованих за допомогою моделі симуляції опилення дикої чорниці (Wild Blueberry Pollination Simulation Model). Ця модель була підтверджена та валідована на основі польових спостережень і експериментальних даних, які були зібрані в штаті Мен (США) та Канадських Марітимах протягом останніх трьох десятиліть.

Подальший аналіз кореляції між різними характеристиками та врожайністю виявив, що деякі ознаки, такі як "fruitset", "seeds" та "fruitmass", мають високу позитивну кореляцію з врожайністю, тоді як інші, наприклад "AverageRainingDays" та "RainingDays", мають виражену негативну кореляцію.

Були виконані два важливі кроки: проведено вибір ключових ознак за допомогою SelectKBest і сформовано ансамбль моделей, базуючись на XGBoost, LightGBM та CatBoost. Вибір ознак спрямований на акцентування найбільш значущих змінних, що сприяє підвищенню якості прогнозування. Ансамбль моделей поєднує можливості трьох сильних алгоритмів машинного навчання, що сприяє покращенню точності прогнозів. У сукупності ці кроки формують робустану основу для інтелектуального модулю прогнозування врожайності ягід.

Інтелектуальний модуль прогнозування, що базується на ансамблі моделей XGBoost, LightGBM та CatBoost, ефективно виконав поставлене завдання. Ці моделі були навчені з використанням оптимальних гіперпараметрів, визначених через процес оптимізації Optuna, що сприяло підвищенню їхньої ефективності. Оглядаючи графіки прогнозованих та дійсних значень, з'ясовується, що вони відображають подібні тренди, свідчаючи про високу спроможність моделей адекватно прогнозувати врожайність. Це підкреслює, що застосування машинного навчання та оптимізації гіперпараметрів може бути надзвичайно корисним для прогнозування врожайності ягід. Така модель може стати в нагоді фермерам та агрономам для оптимізації процесу вирощування ягід та планування урожаю.

					ДП.КН. 8351441.062ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		56

СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ

1. Asare, E., Hoshide, A. K., Drummond, F. A., Criner, G. K., & Chen, X. (2017). Economic risk of bee pollination in Maine wild blueberry, *Vaccinium angustifolium*. *Journal of Economic Entomology*, 110(5), 1980–1992. doi: 10.1093/jee/tox191
2. Popp, J., Pető, K., & Nagy, J. (2013). Pesticide productivity and food security: A review. *Agronomy for Sustainable Development*, 33(1), 243–255. doi: 10.1007/s13593-012-0105-x
3. Drummond, F. (2019). Reproductive biology of wild blueberry (*Vaccinium angustifolium* Aiton). *Agriculture*, 9(4), 69. doi: 10.3390/agriculture9040069
4. Tasnim, R., Drummond, F., & Zhang, Y.-J. (2021). Climate change patterns of wild blueberry fields in downeast, Maine over the past 40 years. *Water*, 13(5), 594. doi: 10.3390/w13050594
5. Tuell, J. K., & Isaacs, R. (2010). Weather during Bloom affects pollination and yield of highbush blueberry. *Journal of Economic Entomology*, 103, 557–562. doi: 10.1603/EC09387
6. Drummond, F. A., Lund, J., & Eitzer, B. (2021). Honey bee health in Maine wild blueberry production. *Insects*, 12(6), 523. doi: 10.3390/insects12060523
7. Crane-Droesch, A. (2018). Machine learning methods for crop yield prediction and climate change impact assessment in agriculture. *Environmental Research Letters*, 13(11), 114003. doi: 10.1088/1748-9326/aae159
8. Ji, B., Sun, Y., Yang, S., & Wan, J. (2007). Artificial neural networks for rice yield prediction in mountainous regions. *Journal of Agricultural Science*, 145(3), 249–261. doi: 10.1017/S0021859606006691
9. Matsumura, K., Gaitan, C. F., Sugimoto, K., Cannon, A. J., & Hsieh, W. W. (2015). Maize yield forecasting by linear regression and artificial neural networks in Jilin, China. *Journal of Agricultural Science*, 153(3), 399–410. doi: 10.1017/S0021859614000392

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		57

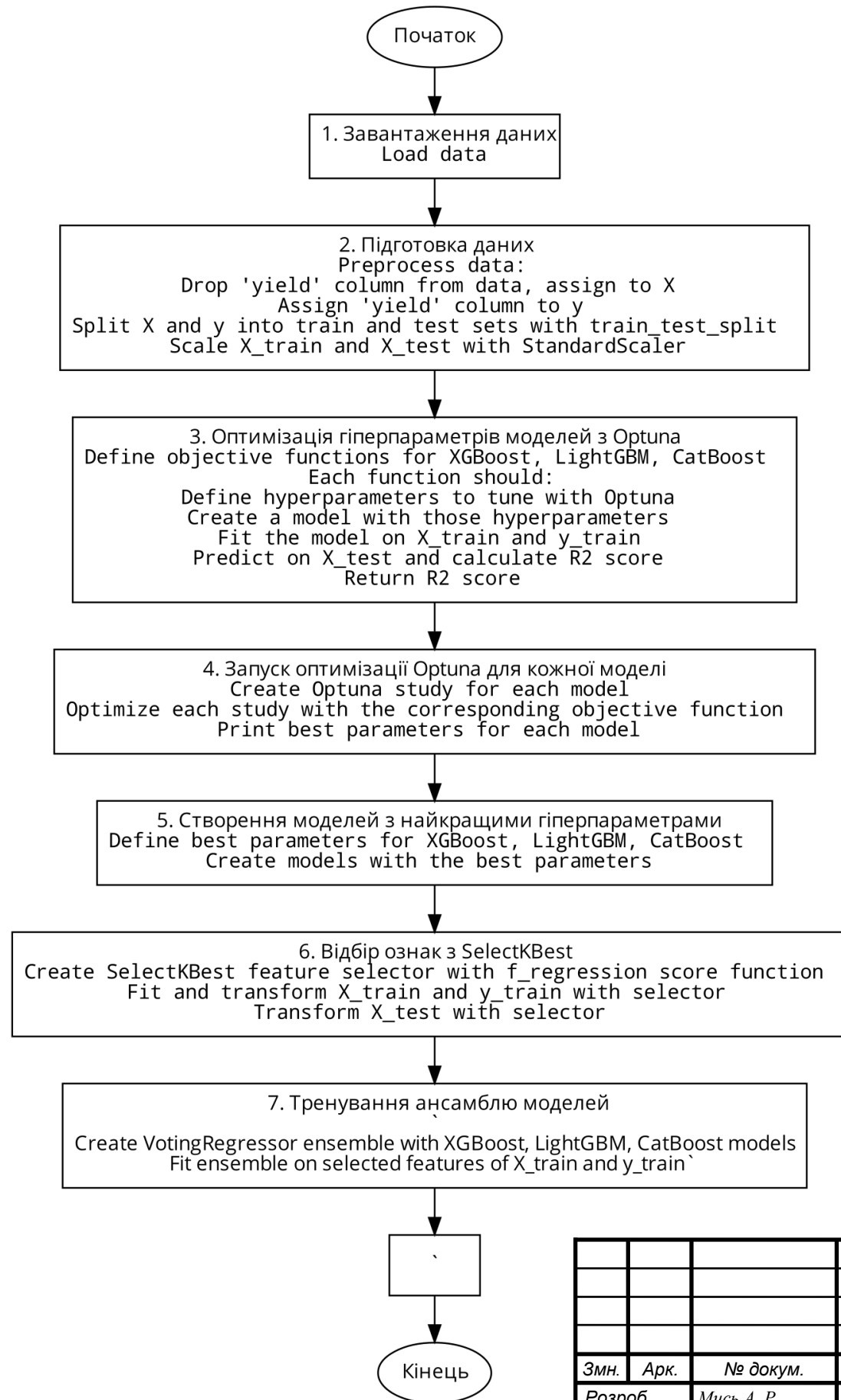
10. Pantazi, X. E., Moshou, D., Alexandridis, T., Whetton, R. L., & Mouazen, A. M. (2016). Wheat yield prediction using machine learning and advanced sensing techniques. *Computers and Electronics in Agriculture*, 121, 57–65. doi: 10.1016/j.compag.2015.11.018
11. Conțiu, Ș., & Groza, A. (2016). Improving remote sensing crop classification by argumentation-based conflict resolution in ensemble learning. *Expert Systems with Applications*, 64, 269–286. doi: 10.1016/j.eswa.2016.08.006
12. Zhang, Y., & Wang, X. (2017). Yield prediction of summer maize in China using remote sensing data and a process-based model. *International Journal of Remote Sensing*, 38(4), 1025–1041. doi: 10.1080/01431161.2016.1268963
13. Balafoutis, A., Beck, P. S. A., Bellon-Maurel, V., & Vangeyte, J. (2020). Machine learning in agriculture: A review. *Sensors*, 20(10), 2881. doi: 10.3390/s20102881
14. Jha, P., Joshi, D., & Singh, P. K. (2019). Machine learning techniques for crop yield prediction: A systematic review. *Computers and Electronics in Agriculture*, 161, 272–293. doi: 10.1016/j.compag.2019.03.014
15. Mohanty, S. P., Hughes, D. P., & Salathé, M. (2016). Using deep learning for image-based plant disease detection. *Frontiers in Plant Science*, 7, 1419. doi: 10.3389/fpls.2016.01419
16. Sladojevic, S., Arsenovic, M., Anderla, A., Culibrk, D., & Stefanovic, D. (2016). Deep neural networks based recognition of plant diseases by leaf image classification. *Computational Intelligence and Neuroscience*, 2016, 3289801. doi: 10.1155/2016/3289801
17. Kamilaris, A., & Prenafeta-Boldú, F. X. (2018). Deep learning in agriculture: A survey. *Computers and Electronics in Agriculture*, 147, 70–90. doi: 10.1016/j.compag.2018.02.016
18. Singh, A., Ganapathysubramanian, B., & Singh, A. K. (2016). Machine learning for high-throughput stress phenotyping in plants. *Trends in Plant Science*, 21(2), 110–124. doi: 10.1016/j.tplants.2015.10.014
19. Mohanty, S. P., Mohanty, S., Hughes, D. P., & Salathé, M. (2017). Using machine learning algorithms to analyze large-scale plant phenomics data generated by

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		58

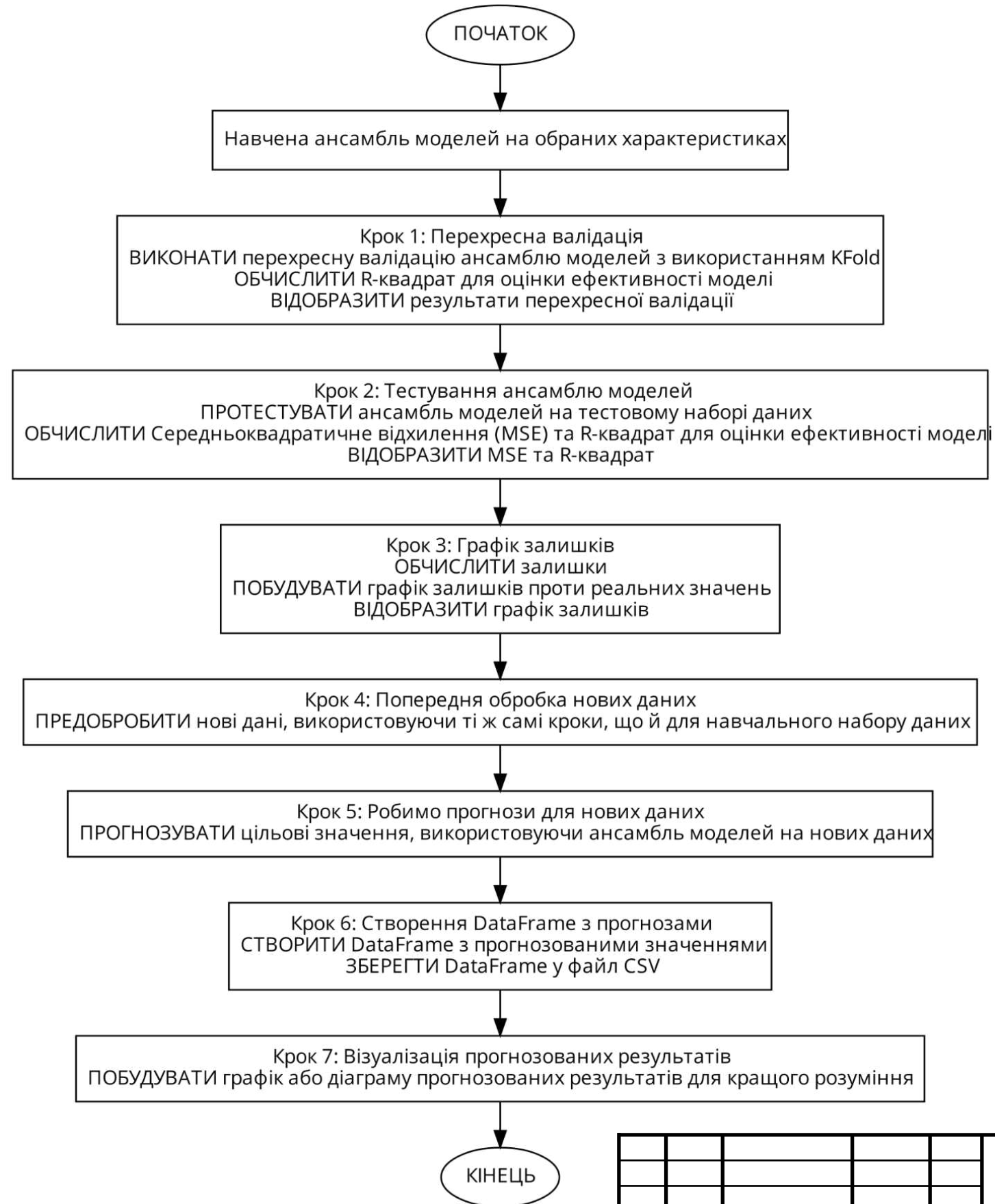
16th International Conference on Computer Sciences and Information Technologies (CSIT) (Vol. 2, pp. 210-216). IEEE.

45. Lipianina-Honcharenko, K., Lukasevych-Krutnyk, I., Butryn-Boka, N., Sachenko, A., & Grodskyi, S. (2021, October). Intelligent Method for Identifying the Fraudulent Online Stores. In *2021 IEEE 8th International Conference on Problems of Infocommunications, Science and Technology (PIC S&T)* (pp. 218-222). IEEE.
46. Komar, M., Savenko, O., Sachenko, A., Lendiuk, T., Lipianina-Honcharenko, K., Hladiy, G., & Vasylkiv, N. (2022). Evaluation the Efficiency of Information Technology of Big Data Intelligence Analysis and Processing.
47. Qu, H. (2020, 12 вересня). Data for: Wild blueberry yield prediction using a combination of computer simulation and machine learning algorithms. Mendeley Data. <https://data.mendeley.com/datasets/p5hvjzsvn8/1>
48. Комар М.П., Саченко А.О., Васильків Н.М., Гладій Г.М., Коваль В.С. Методичні рекомендації до виконання кваліфікаційної роботи з освітньо-професійної програми «Комп'ютерні науки» спеціальності 122 «Комп'ютерні науки» за першим (бакалаврським) рівнем вищої освіти. – Тернопіль: ЗУНУ, 2021. – 56 с.

					<i>ДП.КН. 8351441.062ПЗ</i>	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		62



					ДП.КН. 8351441.001 А1		
					Алгоритм навчання моделей прогнозування врожайності ягід		
					Літера	Маса	Масштаб
Змн.	Арк.	№ докум.	Підпис	Дата			
Розроб.		Мись А. Р.					
Перевір.		Лендюк Т.В.					
Консультант							
Т. Контр.							
Н. Контр.		Лендюк Т.В.					
Затверд.		Комар М.П.					
					Аркуш 1	Аркушів 1	
					ЗУНУ.ФКІТ.КН-41		



					ДП.КН. 8351441.001 А1			
					Алгоритм валідації моделей прогнозування врожайності ягід	Літера	Маса	Масштаб
Змн.	Арк.	№ докум.	Підпис	Дата				
Розроб.		Мись А. Р.						
Перевір.		Лендюк Т.В.						
Консультант								
Т. Контр.						Аркуш 1	Аркушів 1	
Н. Контр.		Лендюк Т.В.				ЗУНУ.ФКІТ.КН-41		
Затверд.		Комар М.П.						