

В. Д. САВЕНКО

Херсонський державний університет

О. Р. ЯРЕМА

Львівський національний університет імені Івана Франка

С. А. БАБІЧЕВ

Університет Яна Евангелиста Пуркіне в Усті на Лабі

Херсонський державний університет

СТРАТЕГІЇ АНСАМБЛЕВОГО НАВЧАННЯ ДЛЯ ПРОГНОЗУВАННЯ УСПІШНОСТІ СТУДЕНТІВ: ПОРІВНЯЛЬНА ТА ІНТЕРАКТИВНА СТРУКТУРА

Ансамблеві методи машинного навчання є одним із найбільш ефективних напрямів побудови інтелектуальних систем аналізу даних, оскільки забезпечують підвищення точності, стійкості та узагальнювальної здатності моделей у задачах класифікації та регресії. Актуальність дослідження зумовлена тим, що сучасні прикладні набори даних часто характеризуються високою розмірністю ознак, наявністю шумів, пропущених значень, викидів і дисбалансу класів, що ускладнює застосування окремих алгоритмів машинного навчання та знижує надійність прогнозування. У статті досліджено теоретичні та практичні засади застосування ансамблевих підходів, зокрема *bagging*, *boosting*, *voting* і *stacking*, з акцентом на їхню роль у підвищенні якості моделювання та зменшенні ризику перенавчання. Особливу увагу приділено стекінгу як гнучкій стратегії інтеграції різнорідних базових моделей за допомогою метамоделі. Практичну частину роботи присвячено розробці інтерактивного веб-інтерфейсу для дослідження ансамблевих методів на прикладі задачі оцінки успішності студентів. У межах запропонованого підходу реалізовано дві стратегії моделювання: пряму багатокласову класифікацію на дискретизованій цільовій змінній та регресію з подальшою категоризацією прогнозованих значень. Для побудови моделей використано алгоритми *Random Forest*, *XGBoost*, *Ridge Regression*, а також ансамблеві схеми *Voting* і *Stacking*. Якість класифікаційних моделей оцінювалася за метриками *Accuracy*, *Precision*, *Recall*, *F1-score* та *Balanced Accuracy*, а регресійних – за показниками *MAE*, *MSE*, *RMSE* і R^2 . У роботі також враховано етапи попередньої обробки даних, крос-валідації та оптимізації гіперпараметрів, що дозволило підвищити відтворюваність і надійність результатів. Розроблений веб-інтерфейс забезпечує поетапне завантаження даних, налаштування параметрів обробки, навчання моделей, аналіз метрик і візуалізацію результатів, що сприяє прозорості експериментування та зручності порівняння різних ансамблевих стратегій. Запропонований підхід є практично придатним для освітніх і дослідницьких задач та може бути використаний як інструмент підтримки прийняття рішень у задачах прогнозування на основі табличних даних.

Ключові слова: Машинне навчання, Ансамблеві методи, Класифікація, Регресія, *Stacking*, *Voting*, *Bagging*, *Boosting*, *Random Forest*, *XGBoost*, Аналіз даних.

V. D. SAVENKO

Kherson State University

O. R. YAREMA

Ivan Franko National University of Lviv

S. A. BABICHEV

Evangelista Purkyně University in Ústí nad Labem

Kherson State University

ENSEMBLE LEARNING STRATEGIES FOR STUDENT PERFORMANCE PREDICTION: A COMPARATIVE AND INTERACTIVE FRAMEWORK

Ensemble methods in machine learning are among the most effective approaches for building intelligent data analysis systems, as they improve the accuracy, robustness, and generalization ability of models in classification and regression tasks. The relevance of this research is due to the fact that modern applied datasets are often characterized by high feature dimensionality, noise, missing values, outliers, and class imbalance, which complicate the use of individual machine learning algorithms and reduce the reliability of predictions. The article examines the theoretical and practical foundations of ensemble approaches, in particular *bagging*, *boosting*, *voting*, and *stacking*, with an emphasis on their role in improving modeling quality and reducing the risk of overfitting. Special attention is paid to *stacking* as a flexible strategy for integrating heterogeneous base models by means of a meta-model. The practical part of the work is devoted to the development of an interactive web interface for studying ensemble methods using the task of student performance

assessment as an example. Within the proposed approach, two modeling strategies are implemented: direct multiclass classification on a discretized target variable and regression followed by categorization of the predicted values. Random Forest, XGBoost, Ridge Regression, as well as Voting and Stacking ensemble schemes were used to build the models. The quality of classification models was evaluated using Accuracy, Precision, Recall, F1-score, and Balanced Accuracy metrics, while regression models were assessed using MAE, MSE, RMSE, and R^2 . The work also takes into account the stages of data preprocessing, cross-validation, and hyperparameter optimization, which made it possible to improve the reproducibility and reliability of the results. The developed web interface provides step-by-step data loading, configuration of preprocessing parameters, model training, metric analysis, and result visualization, which contributes to the transparency of experimentation and the convenience of comparing different ensemble strategies. The proposed approach is practically suitable for educational and research tasks and can be used as a decision-support tool in forecasting problems based on tabular data.

Keywords: Machine learning, Ensemble methods, Classification, Regression, Stacking, Voting, Bagging, Boosting, Random Forest, XGBoost, Data analysis.

Постановка проблеми

Сучасний розвиток методів аналізу даних супроводжується зростанням потреби у точних, стійких та адаптивних моделях прогнозування для задач класифікації та регресії [1; 2]. За умов високої розмірності ознак, наявності шуму, пропущених значень і нерівномірного розподілу даних застосування окремих алгоритмів машинного навчання часто виявляється недостатньо ефективним [1].

Одним із перспективних напрямів підвищення якості моделювання є використання ансамблевих методів, зокрема bagging, boosting, voting та stacking, які дають змогу поєднувати переваги різних моделей і підвищувати узагальнювальну здатність систем [3–6]. Водночас ефективність ансамблевих підходів значною мірою визначається вибором базових моделей, способами їх комбінування та стратегією формування цільової змінної, що зумовлює потребу в системному порівняльному аналізі [3; 5; 6]. Особливої актуальності набуває дослідження різних стратегій моделювання, зокрема підходів прямої багатокласової класифікації та регресії з подальшою категоризацією результатів, що дає змогу гнучкіше враховувати специфіку прикладних задач прогнозування.

Таким чином, актуальною є наукова проблема розробки та порівняльного аналізу підходу до побудови ансамблевих моделей машинного навчання, який забезпечував би можливість інтерактивного дослідження, налаштування та оцінювання ефективності різних ансамблевих стратегій на реальних даних [3; 5; 6].

Аналіз останніх досліджень та публікацій

Проблематика застосування методів машинного навчання у задачах класифікації та регресії активно розвивається як у теоретичному, так і в прикладному аспектах. У фундаментальних працях зі статистичного навчання розглянуто базові принципи побудови прогнозних моделей, зокрема проблеми перенавчання, компромісу між зміщенням і дисперсією, а також роль процедур валідації та регуляризації у підвищенні якості узагальнення [1; 2].

Значна увага в сучасних дослідженнях приділяється ансамблевим методам, які розглядаються як ефективний засіб підвищення точності, стійкості та надійності прогнозних моделей. Показано, що комбінування кількох базових алгоритмів дає змогу зменшити вплив випадкових похибок і покращити результати прогнозування порівняно з використанням окремих моделей [3–6]. Класичні підходи, такі як bagging і його розвиток у вигляді випадкових лісів, продемонстрували високу ефективність для широкого класу задач аналізу даних [7; 8].

Окремий напрям досліджень пов'язаний із розвитком методів стекінгу, які передбачають побудову багаторівневих ансамблевих моделей із використанням метамоделі для інтеграції результатів базових алгоритмів. У працях [9–10] показано, що такий підхід є особливо ефективним у разі використання різнорідних моделей, оскільки дає змогу враховувати їхні індивідуальні переваги та компенсувати окремі недоліки.

У прикладних дослідженнях ансамблеве навчання широко застосовується для розв'язання задач прогнозування в різних предметних областях, зокрема в медицині, фінансовій аналітиці та освітніх системах, де дані характеризуються складною структурою, наявністю шуму та неоднорідністю [11; 12].

Водночас аналіз сучасних публікацій свідчить, що, попри значну кількість праць, недостатньо дослідженими залишаються питання системного порівняння різних ансамблевих стратегій у межах єдиного підходу, зокрема з урахуванням альтернативних способів формування цільової змінної – прямої класифікації та регресії з подальшою категоризацією. Крім того, обмежено представлено інструментальні рішення, які забезпечують можливість інтерактивного дослідження, налаштування та візуального аналізу ансамблевих моделей у прикладних задачах.

Це зумовлює необхідність розробки підходу до побудови та порівняльного аналізу ансамблевих моделей, який поєднує методологічну складову з можливістю інтерактивного експериментування та оцінювання ефективності моделей на реальних даних [5; 6; 9].

Мета дослідження

Мета дослідження полягає у розробці та дослідженні стратегій ансамблевого навчання для задач прогнозування, що включають класифікацію та регресію, зокрема шляхом порівняльного аналізу різних підходів до побудови моделей і формування цільової змінної, а також у практичній реалізації інтерактивного середовища для налаштування, дослідження та оцінювання ефективності ансамблевих моделей на реальних табличних даних.

Виклад основного матеріалу дослідження

Узагальнену послідовність етапів побудови ансамблевої системи машинного навчання, що охоплює підготовку вхідних даних, навчання моделей та оцінювання їх якості, наведено на рис. 1.

На початковому етапі розглядається набір даних, що містить множину вхідних ознак і цільову змінну, яка визначає постановку задачі. Структура даних може включати як числові, так і категоріальні змінні. Водночас реальні набори даних, як правило, характеризуються наявністю пропущених значень, викидів, дублікатів і неоднорідністю форматів подання, що обумовлює необхідність їх попередньої обробки.

Етап попередньої обробки спрямований на забезпечення коректності подальшого моделювання. У межах цього етапу виконуються процедури очищення даних, зокрема обробка пропущених значень і видалення дублікатів, а також трансформація ознак. Остання включає масштабування числових змінних (нормалізація або стандартизація) та кодування категоріальних ознак (зокрема, One-Hot Encoding або Label Encoding). Зазначені перетворення забезпечують приведення даних до формату, придатного для застосування алгоритмів машинного навчання.

Після завершення попередньої обробки визначається тип цільової змінної відповідно до постановки задачі. У задачах класифікації вона задається як дискретна множина класів, тоді як у задачах регресії – як неперервна числова величина. У деяких випадках доцільним є перетворення неперервної змінної на дискретну шляхом її категоризації, що дозволяє перейти до постановки задачі багатокласової класифікації.

На наступному етапі здійснюється навчання базових моделей першого рівня на підготовленому наборі даних. Як базові алгоритми можуть використовуватися різні методи машинного навчання, зокрема Random Forest, XGBoost, Ridge Regression та Ridge Classifier. Вибір моделей визначається специфікою задачі та структурою даних. З метою підвищення узагальнювальної здатності моделей і зменшення ризику перенавчання застосовуються процедури оптимізації гіперпараметрів, зокрема байєсівська оптимізація у поєднанні з k-fold крос-валідацією.

Після побудови базових моделей виконується їх інтеграція в ансамблеву структуру. Найбільш поширеними підходами є voting, averaging та stacking. У випадку використання стекінгу

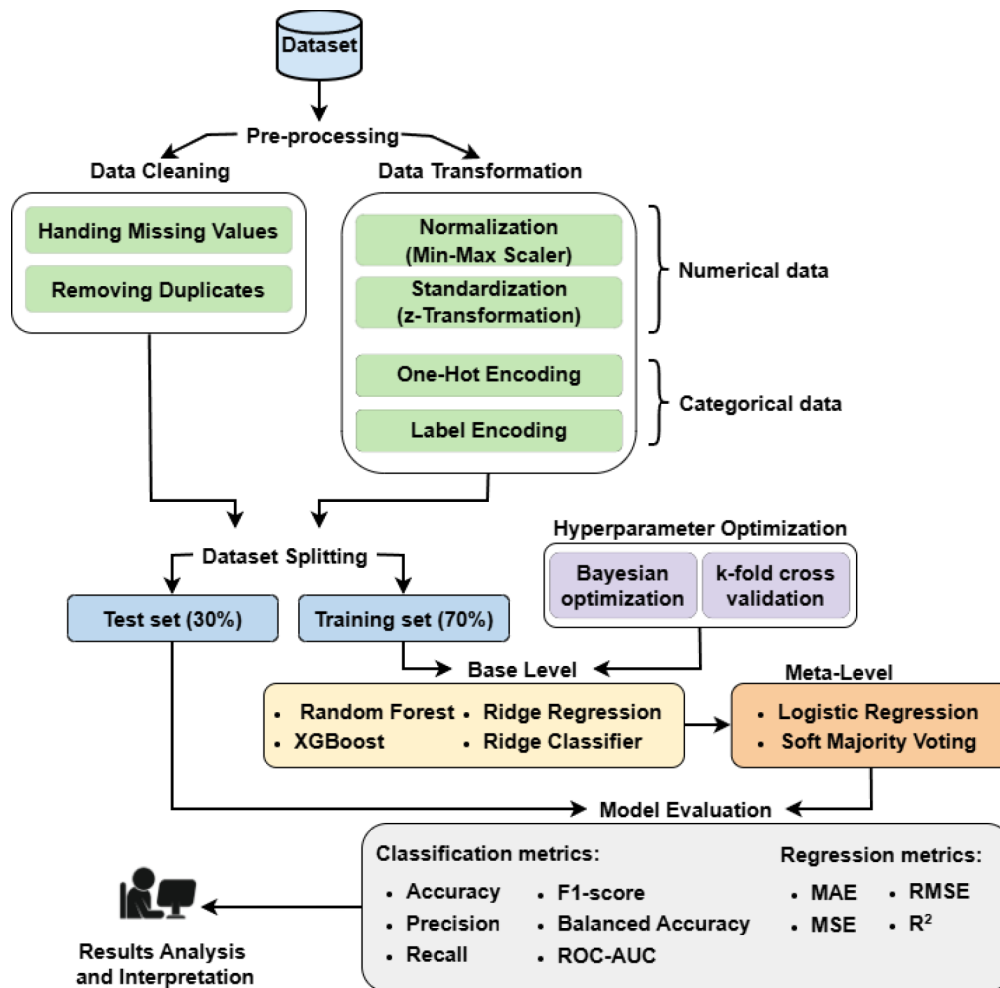


Рис. 1. Структурна схема процесу моделювання

прогнози базових моделей формують простір метаознак, який використовується для навчання метамоделі другого рівня. Такий підхід дозволяє підвищити точність прогнозування за рахунок комбінування переваг окремих алгоритмів і зменшення впливу їхніх індивідуальних похибок.

Завершальним етапом є оцінювання якості побудованих моделей на тестовій вибірці. Вибір метрик визначається типом задачі. Для задач класифікації застосовуються показники асигура, precision, recall, F1-score, balanced accuracy та ROC-AUC. Для задач регресії використовуються метрики MAE, MSE, RMSE та коефіцієнт детермінації R^2 . Аналіз значень зазначених метрик дозволяє здійснити порівняльну оцінку ефективності окремих моделей і ансамблевих підходів.

У межах дослідження розглядається задача прогнозування, яка залежно від природи цільової змінної формулюється як задача класифікації або регресії. Формалізація задачі ґрунтується на заданні навчальної вибірки, визначенні відображення між простором ознак і простором відповідей, а також на виборі функції втрат, мінімізація якої забезпечує навчання моделі.

У задачі класифікації цільова змінна набуває дискретних значень із множини класів:

$$C = \{c_1, c_2, \dots, c_K\}, \tag{1}$$

де K – кількість класів.

У задачі регресії цільова змінна є неперервною величиною:

$$y \in \mathbb{R}. \tag{2}$$

Навчальна вибірка задається у вигляді множини пар «ознаки–відповідь»:

$$D = \{(x_i, y_i)\}_{i=1}^N, \tag{3}$$

де $x_i = (x_{i1}, x_{i2}, \dots, x_{im})$ – вектор ознак i -го об’єкта, $y_i \in Y$ – відповідне значення цільової змінної, N – кількість об’єктів у вибірці, m – кількість ознак.

У загальному випадку задача навчання з учителем зводиться до знаходження відображення:

$$f: X \rightarrow Y, \quad (4)$$

де $X \in \mathbb{R}^m$ – простір вхідних ознак, Y – простір значень цільової змінної.

Для задачі класифікації відображення маємо $f: X \rightarrow \mathcal{C}$, а для задачі регресії – $f: X \rightarrow \mathbb{R}$. Якість моделі оцінюється за допомогою функції втрат $L(y_i, f(x_i))$, яка визначає ступінь невідповідності між істинним значенням цільової змінної та прогнозом моделі. На навчальній вибірці вводиться емпіричний ризик:

$$Q(f) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i)). \quad (5)$$

Метою навчання є знаходження оптимального відображення:

$$f^* = \arg \min_f Q(f). \quad (6)$$

У задачах класифікації як функції втрат застосовуються, зокрема, індикаторна функція помилки або логарифмічна втрата, тоді як у задачах регресії найчастіше використовуються середньоквадратична та середньоабсолютна похибки.

Таким чином, задачі класифікації та регресії відрізняються природою простору відповідей, однак мають спільну математичну основу, що полягає у побудові відображення між простором ознак і простором відповідей на основі навчальної вибірки шляхом мінімізації функції втрат.

У межах дослідження розглянуто два підходи до формування цільової змінної в задачі прогнозування академічної успішності. Запропоновані підходи відрізняються способом задання цільової змінної та механізмом переходу до кінцевих класів і розглядаються в межах єдиної дослідницької схеми, що забезпечує можливість їх коректного порівняльного аналізу.

Перший підхід ґрунтується на прямій класифікації, за якої цільова змінна дискретизується до початку навчання моделі. Нехай початкове значення цільової змінної є неперервним: $y \in \mathbb{R}$. Тоді дискретна цільова змінна визначається як:

$$y^{(c)} = \phi(y), \quad (7)$$

де $\phi(\cdot)$ – функція дискретизації, що відображає неперервне значення y у відповідний клас $y^{(c)} \in \mathcal{C} = \{c_1, c_2, \dots, c_K\}$.

У цьому випадку навчальна вибірка набуває вигляду

$$\mathcal{D}_{cls} = \{(x_i, y_i^{(c)})\}_{i=1}^N, \quad (8)$$

а модель безпосередньо реалізує відображення

$$f_{cls}: \mathcal{X} \rightarrow \mathcal{C}. \quad (9)$$

Таким чином, результатом роботи моделі є одразу прогноз класової належності об’єкта без проміжного етапу оцінювання числового значення.

Другий підхід передбачає схему регресії з подальшою категоризацією, у межах якої модель спочатку прогнозує неперервне значення цільової змінної, а вже потім це значення перетворюється на клас. Спочатку будується регресійна модель

$$f_{reg}: \mathcal{X} \rightarrow \mathbb{R}, \quad (10)$$

яка для кожного об’єкта формує прогноз

$$\hat{y} = f_{reg}(x). \quad (11)$$

Далі, прогнозоване неперервне значення перетворюється на клас за допомогою функції категоризації

$$\hat{y}^{(c)} = \phi(\hat{y}), \quad (12)$$

де $\phi(\cdot)$ задається системою порогів і відображає прогнозоване числове значення у відповідний клас $\hat{y}^{(c)}$.

Принципова відмінність між розглянутими підходами полягає у характері оптимізаційної задачі. У випадку прямої класифікації модель безпосередньо мінімізує класифікаційну функцію втрат відносно дискретної цільової змінної. Натомість у підході регресії з подальшою категоризацією мінімізується похибка прогнозування неперервної змінної, тоді як віднесення до класів виконується на постобробному етапі.

Таким чином, у першому випадку межі між класами враховуються безпосередньо під час навчання, тоді як у другому – опосередковано, через застосування функції дискретизації до прогнозованих значень.

Запропоновані стратегії формування цільової змінної дозволяють розглядати задачу прогнозування в альтернативних постановках у межах єдиного підходу, що забезпечує можливість їх системного порівняння за показниками точності, стійкості та інтерпретованості.

Як базові моделі ансамблю використано алгоритми Random Forest, XGBoost та Ridge. Вибір зазначених моделей обумовлено необхідністю забезпечення різноманітності прогнозів, що є ключовим чинником ефективності ансамблевих підходів. Обрані алгоритми представляють різні парадигми машинного навчання. Зокрема, Random Forest реалізує підхід bagging і спрямований на зменшення дисперсії прогнозу шляхом агрегування результатів множини дерев рішень. XGBoost належить до методів boosting і забезпечує послідовне уточнення моделі за рахунок мінімізації похибок попередніх ітерацій. Ridge є лінійною моделлю з L_2 -регуляризацією, що забезпечує стабільність оцінок параметрів і характеризується високою інтерпретованістю.

Таким чином, поєднання моделей, що реалізують різні підходи до побудови прогнозу (bagging, boosting та регуляризовані лінійні моделі), дозволяє враховувати як нелінійні, так і лінійні залежності у даних, а також підвищує стійкість і узагальнювальну здатність ансамблевої системи.

Для поєднання прогнозів базових моделей використано два основні підходи ансамблювання: Voting та Stacking. Їх застосування дає змогу об'єднати результати моделей різної природи та підвищити стійкість і точність кінцевого прогнозу. Якщо метод voting реалізує безпосереднє агрегування прогнозів базових алгоритмів, то stacking передбачає побудову дворівневої структури, у якій фінальне рішення формується окремою метамоделлю. Саме тому в межах дослідження stacking розглядається як центральний елемент ансамблювання, оскільки він дозволяє не лише поєднати прогнози, а й навчити модель другого рівня оптимально використовувати інформацію, отриману від базових алгоритмів.

У методі Voting фінальний прогноз формується шляхом агрегування виходів базових моделей. Для задачі класифікації це може бути правило більшості голосів, яке формально задається співвідношенням

$$\hat{y} = \arg \max_{c \in C} \sum_{j=1}^M \mathbb{I}(f_j(x) = c), \quad (13)$$

де M – кількість базових моделей, $f_j(x)$ – прогноз j -ї моделі, C – множина класів, $\mathbb{I}(\cdot)$ – індикаторна функція.

У разі використання зваженого або soft voting агрегування може здійснюватися на основі прогнозованих імовірностей:

$$\hat{y} = \arg \max_{c \in C} \sum_{j=1}^M w_j p_j(cx), \quad (14)$$

де w_j – вага j -ї моделі, $p_j(c|x)$ – оцінена ймовірність належності об'єкта до класу c .
Для задачі регресії агрегування прогнозів зазвичай виконується усередненням

$$\hat{y} = \frac{1}{M} \sum_{j=1}^M f_j(x). \quad (15)$$

На відміну від voting, метод stacking реалізує ієрархічну схему ансамблювання. На першому рівні базові моделі формують індивідуальні прогнози

$$z_j = f_j(x), \quad j = 1, 2, \dots, M, \quad (16)$$

які використовуються як метаознаки для моделі другого рівня. У такому разі вектор метаознак має вигляд

$$z = (z_1, z_2, \dots, z_M), \quad (17)$$

а фінальний прогноз визначається метамоделлю

$$\hat{y} = g(z) = g(f_1(x), f_2(x), \dots, f_M(x)), \quad (18)$$

де $g(\cdot)$ – метамоделль.

Отже, роль метамоделі полягає в тому, щоб навчитися враховувати сильні сторони базових алгоритмів, компенсувати їх слабкі місця та формувати більш точний підсумковий прогноз, ніж просте усереднення або голосування.

Таким чином, voting забезпечує просте та інтерпретоване агрегування прогнозів, тоді як stacking реалізує більш гнучкий механізм ансамблювання за рахунок використання метамоделі другого рівня. Саме тому stacking доцільно розглядати як основний метод побудови ансамблю в даному дослідженні, оскільки він дозволяє ефективно комбінувати прогнози різнорідних базових моделей і підвищувати якість кінцевого результату.

Для оцінювання якості побудованих моделей використано окремі групи метрик для задач класифікації та регресії. Такий поділ зумовлений різною природою цільової змінної та, відповідно, різними підходами до інтерпретації похибки прогнозування.

Для задач класифікації застосовано метрики Accuracy, Precision, Recall, F1-score та Balanced Accuracy. Метрика Accuracy визначає частку правильно класифікованих об'єктів серед усіх спостережень:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \quad (19)$$

де TP – кількість істинно позитивних прогнозів, TN – кількість істинно негативних прогнозів, FP – кількість хибнопозитивних прогнозів, FN – кількість хибнонегативних прогнозів.

Метрика Precision характеризує частку правильно визначених позитивних об'єктів серед усіх об'єктів, віднесених моделлю до позитивного класу:

$$Precision = \frac{TP}{TP + FP}. \quad (20)$$

Метрика Recall відображає частку правильно виявлених позитивних об'єктів серед усіх істинно позитивних:

$$Recall = \frac{TP}{TP + FN}. \quad (21)$$

F1-score є гармонійним середнім між Precision і Recall:

$$F1 = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}. \quad (22)$$

Метрика Balanced Accuracy використовується у випадках незбалансованих класів і визначається як середнє значення Recall за всіма класами; для бінарної класифікації її можна подати як

$$Balanced Accuracy = \frac{1}{2} \left(\frac{TP}{TP + FN} + \frac{TN}{TN + FP} \right). \quad (23)$$

Для задач регресії використано метрики MAE, MSE та RMSE, які характеризують величину відхилення прогнозованих значень від істинних. MAE (Mean Absolute Error) визначає середню абсолютну похибку прогнозу:

$$MAE = \frac{1}{N} \sum_{i=1}^N |y_i - \hat{y}_i|. \quad (24)$$

MSE (Mean Squared Error) визначає середню квадратичну похибку:

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2. \quad (25)$$

RMSE (Root Mean Squared Error) є квадратним коренем із середньоквадратичної похибки:

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}. \quad (26)$$

Коефіцієнт детермінації (R^2) відображає частку дисперсії залежної змінної, що пояснюється моделлю:

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}. \quad (27)$$

Використаний набір метрик дає змогу комплексно оцінити якість моделей у межах обох стратегій прогнозування: у задачах класифікації – з погляду правильності розпізнавання класів, а в задачах регресії – з погляду точності прогнозування неперервного значення.

Узагальнена математична модель дослідження формалізується як багаторівнева система прогнозування, у якій простір ознак, множина базових моделей, механізм ансамблювання та функція втрат інтегруються в єдину композиційну структуру. Нехай задано навчальну вибірку (3). Множина базових моделей визначається як

$$F = \{f_1, f_2, \dots, f_M\}, \quad (28)$$

де M – кількість базових алгоритмів.

У межах дослідження такими моделями виступають Random Forest, XGBoost та Ridge, кожна з яких формує індивідуальний прогноз. Сукупність прогнозів базових моделей утворює вектор метаознак, який використовується для побудови ансамблевого прогнозу. У загальному вигляді фінальна ансамблева модель може бути подана як композиція (18). Для voting функція $g(\cdot)$ реалізує безпосереднє агрегування прогнозів базових моделей, а для stacking – відображає роботу метамоделі другого рівня, яка отримує вектор метаознак z і формує остаточний прогноз. Таким чином, якщо voting задає фіксовану схему комбінування, то stacking забезпечує адаптивне навчання правила інтеграції базових моделей.

Якість узагальненої системи оцінюється через функцію втрат

$$L(y, \hat{y}) = L(y, G(F, x)), \quad (29)$$

яка визначає невідповідність між істинним і прогнозованим значеннями цільової змінної, а задача навчання формулюється як мінімізація емпіричного ризику

$$Q(G) = \frac{1}{N} \sum_{i=1}^N L(y_i, G(F, x_i)). \quad (30)$$

У такому вигляді модель об'єднує всі основні компоненти дослідження в єдину систему: вхідні ознаки x , множину базових моделей F , механізм ансамблювання G та функцію втрат L . Отже, узагальнена математична модель відображає повний процес побудови прогнозу – від індивідуальних передбачень базових алгоритмів до формування підсумкового ансамблевого рішення та оцінювання його якості.

Запропонована узагальнена математична модель ансамблевого прогнозування була реалізована у вигляді інтерактивного програмного середовища, що забезпечує можливість поетапного виконання всіх етапів побудови моделі – від завантаження та попередньої обробки даних до навчання базових алгоритмів, формування ансамблю та оцінювання якості отриманих результатів.

Розроблений web-інтерфейс відображає структуру запропонованого підходу та реалізує його як послідовність взаємопов'язаних функціональних блоків, що відповідають основним компонентам математичної моделі: простору ознак X , множині базових моделей F , механізму ансамблювання G та функції оцінювання якості L . Інтерфейс надає користувачу можливість:

- завантаження та попередньої обробки даних;
- вибору стратегії формування цільової змінної (класифікація або регресія з подальшою категоризацією);
- налаштування параметрів базових моделей;
- вибору методу ансамблювання (voting або stacking);
- проведення навчання моделей і аналізу отриманих результатів за заданими метриками.

Загальний вигляд розробленого інтерфейсу наведено на рис. 2. З використанням розробленого середовища було проведено експериментальне дослідження ефективності запропонованих стратегій ансамблевого навчання. Основна увага приділялася порівняльному аналізу двох підходів до формування цільової змінної – прямої класифікації та регресії з подальшою категоризацією – а також оцінюванню впливу методів ансамблювання (voting і stacking) на якість прогнозування.

У дослідженні використано набір даних Student Performance, що містить 10 000 спостережень і 6 ознак, які характеризують навчальну діяльність студентів та рівень їхньої академічної успішності. Цільовою змінною в обох стратегіях був показник Performance Index. Для забезпечення коректного міжстратегічного порівняння розподіл на класи здійснювався за однаковими порогоми $Q_1 = 40$ та $Q_3 = 71$. Формування метаознак виконувалося з використанням 5-кратної крос-валідації, а на другому рівні досліджувалися ансамблі Voting та Stacking.

У межах стратегії прямої багатокласової класифікації було побудовано базові моделі RidgeClassifier, RandomForestClassifier та XGBClassifier, а також ансамблі другого рівня SoftVoting і Stacking. Отримані результати засвідчили, що найвищу якість класифікації забезпечив ансамбль SoftVoting, для якого значення $F1_{weighted}$ становило 0.9487. Це перевищує показники всіх базових моделей першого рівня. Найкращою серед окремих базових моделей виявилася XGBClassifier із $F1_{weighted} = 0.9456$. Для RidgeClassifier і RandomForestClassifier значення $F1_{weighted}$ становили по 0.9440. Порівняння основних метрик якості наведено в табл. 1.

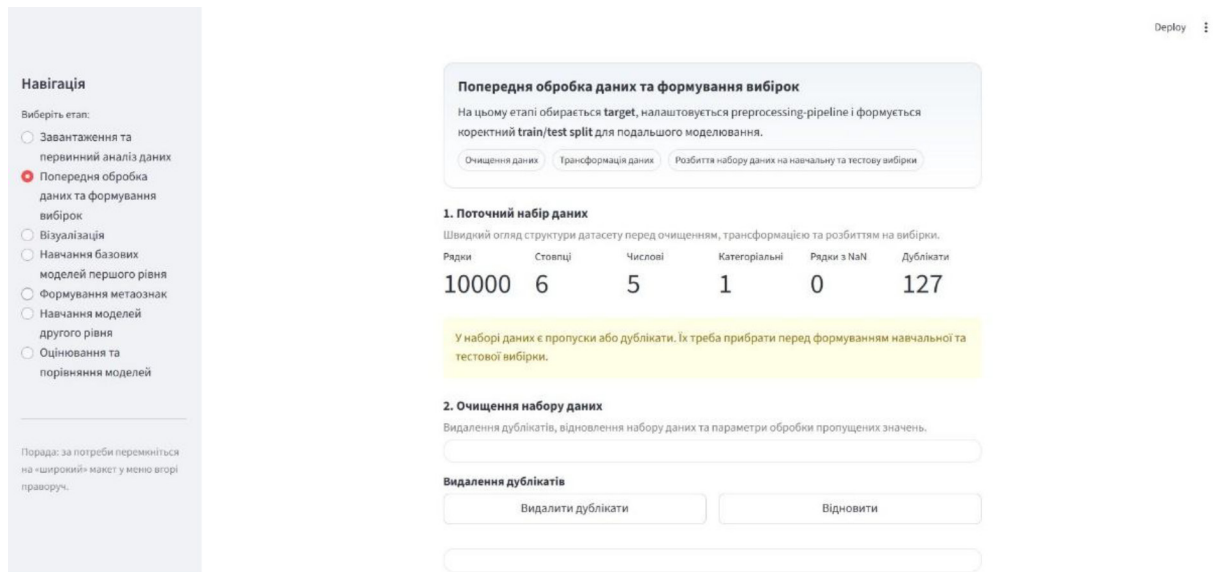


Рис. 2. Загальний вигляд розробленого інтерфейсу

Таблиця 1

Метрики якості моделей у стратегії прямої багатокласової класифікації

Модель	Accuracy	Balanced Accuracy	F1 weighted
Ridge	0.9440	0.9416	0.9440
Random Forest	0.9440	0.9431	0.9440
XGBoost	0.9456	0.9467	0.9456
SoftVoting	0.9487	0.9483	0.9487
Stacking	0.9446	0.9478	0.9446

Перевага SoftVoting у цій стратегії пояснюється використанням імовірностей належності до класів, сформованих базовими моделями, що забезпечує більш точне врахування невизначеності прогнозу та покращує якість фінального рішення.

У межах стратегії регресії з подальшою категоризацією як базові моделі використано RidgeRegressor, RandomForestRegressor та XGBRegressor, а на другому рівні – ансамблі Voting і Stacking. У цій постановці моделі спочатку прогнозували неперервне значення Performance Index, після чого отримані прогнози категоризувалися за тими самими порогами Q_1 та Q_3 . Результати оцінювання наведено в табл. 2.

Таблиця 2

Метрики якості моделей у стратегії регресії з подальшою категоризацією

Модель	R^2	Accuracy	Balanced Accuracy	F1 weighted
Ridge	0.9884	0.9450	0.9450	0.9449
Random Forest	0.9867	0.9429	0.9431	0.9429
XGBoost	0.9881	0.9456	0.9444	0.9456
Voting	0.9882	0.9443	0.9441	0.9442
Stacking	0.9884	0.9453	0.9452	0.9452

Найкращий результат у межах цієї стратегії продемонстрував XGBRegressor після категоризації прогнозів: значення $F1_{weighted}$ становило 0.9456. Ансамблеві моделі Voting і Stacking не забезпечили переваги над найкращою базовою моделлю, хоча їхні результати залишилися близькими. Матриці плутанини для найкращих моделей двох стратегій наведено на рис. 3.

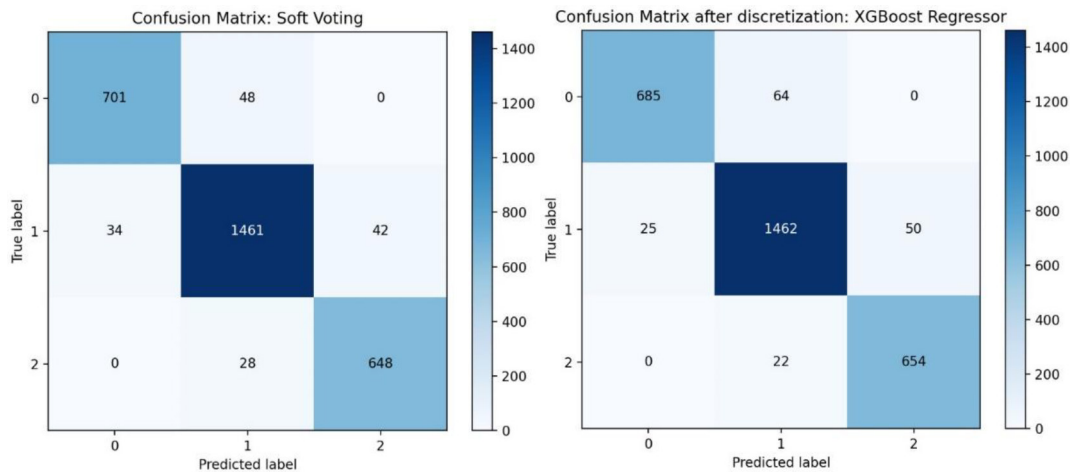


Рис. 3. Матриці плутанини ансамблю SoftVoting та моделі XGBRegressor

Порівняльний аналіз показав, що стратегія прямої багатокласової класифікації виявилася дещо ефективнішою за головною метрикою якості. Найкраща модель цієї стратегії – SoftVoting – досягла $F1_{weighted} = 0.9487$, тоді як у стратегії регресії з подальшою категоризацією найкращий результат становив 0.9456 і був отриманий для XGBRegressor після дискретизації прогнозів. Отримана різниця є незначною, однак свідчить про доцільність безпосереднього розв’язання задачі класифікації для даного набору даних.

Отже, проведене експериментальне дослідження підтвердило ефективність ансамблевих методів для прогнозування успішності студентів. Для досліджуваного набору даних найкращий результат забезпечила стратегія прямої багатокласової класифікації із застосуванням ансамблю SoftVoting, тоді як у регресійній постановці найвищу якість продемонстрував XGBRegressor. Це підтверджує, що вибір стратегії моделювання та способу агрегування прогнозів істотно впливає на підсумкову якість моделі.

Висновки

У статті досліджено дві стратегії ансамблевого навчання для прогнозування успішності студентів: пряму багатокласову класифікацію та регресію з подальшою категоризацією. Реалізовано інтерактивне програмне середовище, що підтримує завантаження, попередню обробку, візуалізацію даних, навчання базових моделей та ансамблів, а також порівняльне оцінювання результатів.

Експериментально встановлено, що для набору даних Student Performance найкращий результат забезпечує стратегія прямої багатокласової класифікації із застосуванням ансамблю SoftVoting ($F1_{weighted} = 0.9487$). У стратегії регресії з подальшою категоризацією найвищу якість продемонструвала модель XGBRegressor після дискретизації прогнозів ($F1_{weighted} = 0.9456$). Отримані результати підтверджують ефективність ансамблевих підходів і показують, що вибір стратегії моделювання суттєво впливає на підсумкову якість прогнозу.

Практична цінність роботи полягає у створенні web-інтерфейсу для інтерактивного застосування ансамблевих методів машинного навчання. Подальші дослідження доцільно спрямувати на розширення набору моделей, оптимізацію ваг ансамблів та тестування підходу на інших наборах даних.

Список використаної літератури

1. Hastie T., Tibshirani R., Friedman J. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd ed. Springer, 2009. DOI: <https://doi.org/10.1007/978-0-387-84858-7>
2. Ghojogh B., Crowley M. *The Theory Behind Overfitting, Cross Validation, Regularization, Bagging, and Boosting: Tutorial*. arXiv, 2019. URL: [arXiv:1905.12787](https://arxiv.org/abs/1905.12787)
3. Zhou Z.-H. *Ensemble Methods: Foundations and Algorithms*. 2nd ed. Chapman and Hall/CRC, 2025. DOI: <https://doi.org/10.1201/9781003587774>
4. Rokach L. Ensemble-based classifiers. *Artificial Intelligence Review*. 2010. Vol. 33, no. 1–2. P. 1–39. DOI: <https://doi.org/10.1007/s10462-009-9124-7>
5. Dietterich T. G. Ensemble methods in machine learning. In: *Multiple Classifier Systems*. Springer, 2000. P. 1–15. DOI: https://doi.org/10.1007/3-540-45014-9_1
6. Kuncheva L. I. *Combining Pattern Classifiers: Methods and Algorithms*. 2nd ed. John Wiley & Sons, 2014.
7. Breiman L. Bagging predictors. *Machine Learning*. 1996. Vol. 24, no. 2. P. 123–140. DOI: <https://doi.org/10.1023/A:1018054314350>
8. Breiman L. Random forests. *Machine Learning*. 2001. Vol. 45, no. 1. P. 5–32. DOI: <https://doi.org/10.1023/A:1010933404324>
9. Wolpert D. H. Stacked generalization. *Neural Networks*. 1992. Vol. 5, no. 2. P. 241–259. DOI: [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1)
10. Sill J., Takács G., Mackey L., Lin D. *Feature-Weighted Linear Stacking*. arXiv, 2009. URL: [arXiv:0911.0460](https://arxiv.org/abs/0911.0460)
11. Shindo J. H., Mjahidi M. M., Waziri M. D. Data mining algorithms for prediction of student teachers' performance in ICT: A systematic literature review. *Information Technologies and Learning Tools*. 2023. Vol. 96, no. 4. P. 29–45. DOI: <https://doi.org/10.33407/itlt.v96i4.5246>
12. Caprian I. Impact of false alarms in machine learning-based anti-fraud systems: The economic and reputational consequences. *Business Inform.* 2025. No. 8. P. 378–389. DOI: <https://doi.org/10.32983/2222-4459-2025-8-378-389>

References

1. Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.). Springer <https://doi.org/10.1007/978-0-387-84858-7> [in English].
2. Ghojogh, B., & Crowley, M. (2019). *The theory behind overfitting, cross validation, regularization, bagging, and boosting: Tutorial*. arXiv [in English].
3. Zhou, Z.-H. (2025). *Ensemble Methods: Foundations and Algorithms* (2nd ed.). Chapman and Hall/CRC <https://doi.org/10.1201/9781003587774> [in English].
4. Rokach, L. (2010). Ensemble-based classifiers. *Artificial Intelligence Review*, 33 (1–2), 1–39 <https://doi.org/10.1007/s10462-009-9124-7> [in English].
5. Dietterich, T. G. (2000). Ensemble methods in machine learning. In *Multiple Classifier Systems* (pp. 1–15). Springer https://doi.org/10.1007/3-540-45014-9_1 [in English].
6. Kuncheva, L. I. (2014). *Combining Pattern Classifiers: Methods and Algorithms* (2nd ed.). John Wiley & Sons [in English].
7. Breiman, L. (1996). Bagging predictors. *Machine Learning*, 24 (2), 123–140 <https://doi.org/10.1023/A:1018054314350> [in English].
8. Breiman, L. (2001). Random forests. *Machine Learning*, 45 (1), 5–32 <https://doi.org/10.1023/A:1010933404324> [in English].
9. Wolpert, D. H. (1992). Stacked generalization. *Neural Networks*, 5 (2), 241–259 [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1) [in English].

10. Sill, J., Takács, G., Mackey, L., & Lin, D. (2009). *Feature-weighted linear stacking*. arXiv [in English].
11. Shindo, J. H., Mjahidi, M. M., & Waziri, M. D. (2023). Data mining algorithms for prediction of student teachers' performance in ICT: A systematic literature review. *Information Technologies and Learning Tools*, 96 (4), 29–45 <https://doi.org/10.33407/itlt.v96i4.5246> [in English].
12. Caprian, I. (2025). Impact of false alarms in machine learning-based anti-fraud systems: The economic and reputational consequences. *Business Inform*, 8, 378–389 <https://doi.org/10.32983/2222-4459-2025-8-378-389> [in English].

Савенко Валерія Дмитрівна – магістрантка кафедри комп'ютерних наук та програмної інженерії Херсонського державного університету. E-mail: 001014@university.kherson.ua, ORCID: 0009-0004-3352-0835.

Ярема Олег Романович – к.т.н., доцент, доцент кафедри цифрової економіки та бізнес аналітики Львівського національного університету імені Івана Франка. E-mail: oleh.yarema@lnu.edu.ua, ORCID: 0000-0003-3736-4820.

Бабічев Сергій Анатолійович – д.т.н., професор, професор кафедри інформатики університету Яна Євангеліста Пуркіне в Усті на Лабі, Чехія; професор кафедри фізики Херсонського державного університету. E-mail: sergii.babichev@ujep.cz, sbabichev@ksu.ks.ua, ORCID: 0000-0001-6797-1467.

Savenko Valeriia Dmytrivna – Master's Student at the Department of Computer Science and Software Engineering of the Kherson State University. E-mail: 001014@university.kherson.ua, ORCID: 0009-0004-3352-0835.

Yarema Oleh Romanovych – Candidate of Technical Sciences, Associate Professor, Associate Professor at the Department of Digital Economy and Business Analytics of the Ivan Franko National University of Lviv. E-mail: oleh.yarema@lnu.edu.ua, ORCID: 0000-0003-3736-4820.

Babichev Sergii Anatoliiiovych – Doctor of Technical Sciences, Professor, Professor at the Department of Informatics of the Jan Evangelista Purkyně University in Ústí nad Labem, Czech Republic; Professor at the Department of Physics of the Kherson State University. E-mail: sergii.babichev@ujep.cz, sbabichev@ksu.ks.ua, ORCID: 0000-0001-6797-1467.

Дата першого надходження статті до видання: 09.04.2026

Дата прийняття статті до друку після рецензування: 14.05.2026

Дата публікації (оприлюднення) статті: 01.07.2026



Стаття поширюється на умовах ліцензії відкритого доступу (CC BY 4.0)