



О. В. Мокрицька¹, Ю. М. Мочернюк²

¹ Національний університет Львівська політехніка, м. Львів, Україна

² Національний лісотехнічний університет України, м. Львів, Україна

ВИКОРИСТАННЯ АЛГОРИТМІВ МАШИННОГО НАВЧАННЯ ДЛЯ АВТОМАТИЗАЦІЇ ПРОЦЕСУ МОДЕРАЦІЇ КОНТЕНТУ В ГРУПОВИХ ЧАТАХ МЕСЕНДЖЕРІВ

Проаналізовано особливості інтеграції методів машинного навчання ML (англ. *Machine Learning*) в автоматизовану систему модераторів групових чатів Telegram для вдосконалення процесу управління широкомасштабними комунікаціями. З'ясовано, що проблема модераторів великих групових чатів стає особливо нагальною через збільшення обсягу повідомлень і розмаїття контенту, що потребує ефективної системи їх фільтрації для забезпечення безпеки комунікаційного середовища. Застосовано розширені алгоритми вибору функцій класифікації, такі як оптимізація рою частинок PSO (англ. *Particle Swarm Optimization*), алгоритм рою салпів SSA (англ. *Salp Swarm Algorithm*) і оптимізація Сірого вовка GWO (англ. *Gray Wolf Optimization*), для автоматизації процесу виявлення найбільш релевантних характеристик повідомлень, що значно покращує процес модераторів групових чатів. Досліджено особливості застосування методів класифікації, зокрема машин підтримки векторів SVM (англ. *Support Vector Machines*) та алгоритму найближчих сусідів kNN (англ. *k-Nearest Neighbors*), для ідентифікації неприйнятної або шкідливого контенту. Встановлено, що вибір класифікаційних функцій є критичним для зменшення розмірності даних і підвищення точності класифікації. Проведено порівняльний аналіз ефективності алгоритмів PSO, SSA і GWO в контексті модераторів чатів. Показано, що алгоритм PSO демонструє високу ефективність завдяки швидкій адаптації до мінливого комунікаційного середовища, тоді як алгоритми SSA і GWO демонструють варіації у підходах до глобальної оптимізації вибору функцій класифікації. Розроблено систему класифікації, в якій комбінація алгоритмів PSO+SVM досягає найвищої точності, перевершуючи інші алгоритми за збалансованістю швидкості та точності класифікації. Наведено також результати для комбінацій алгоритмів SSA+kNN і алгоритмів GWO+kNN, які, хоча і показують високу ефективність, мають менш стабільні результати в різних умовах. Проведено оцінювання продуктивності зазначених алгоритмів за такими показниками, як точність, чутливість, специфічність, вивірність і оцінка F1. Комплексний аналіз цих алгоритмів підтвердив, що впровадження вдосконалених алгоритмів вибору функцій класифікації значно покращує точність виявлення шкідливого контенту, зменшуючи кількість помилкових спрацювань і підвищуючи загальну ефективність процесу модераторів. Дослідження відповідних алгоритмів щодо модераторів вмісту також встановило потенційне застосування розробленої системи на інших платформах та у нових доменах для забезпечення масштабованого й адаптованого рішення та підтримки безпеки в онлайн-комунікаціях.

Ключові слова: Telegram-бот; оптимізація рою частинок; алгоритм рою салпів; оптимізація Сірого вовка; машини підтримки векторів; вибір класифікаційних функцій; класифікація; ансамблеві методи.

Вступ / Introduction

Ефективна модераторів повідомлень та вмісту на широкомасштабних комунікаційних платформах, таких як Telegram, стає дедалі складнішою проблемою в умовах постійного зростання кількості користувачів і повідомлень. Групові чати є місцем інтенсивної взаємодії, однак їхнє модераторів середовище стикається з труднощами, такими як спам, неприйнятний контент і загрози безпеці. Традиційні системи модераторів, засновані на фіксованих правилах, виявляються недостатньо ефективними через неможливість адаптуватися до нових видів зловмисної поведінки та великих обсягів даних. За даними дослідження [9], подібні системи мають низьку точність і високу кількість хибнопозитивних спрацювань.

Машинне навчання останніми роками активно досліджують як можливе рішення для автоматизації процесів модераторів контенту. Автори дослідження [12] запропонували використання моделей класифікації, зокрема Support Vector Machines (SVM) та k-Nearest Neighbors (kNN), для виявлення небезпечного контенту. Однак, існує потреба в удосконаленні процесу вибору функцій класифікації для покращення точності та зменшення кількості помилкових спрацювань. Попри значні досягнення, такі особливості, як вибір найбільш релевантних характеристик повідомлень, залишаються недостатньо дослідженими. Актуальність цього дослідження полягає у потребі підвищення точності та адаптивності модераторів, що забезпечить безпечніше середовище в умовах масштабних онлайн-комунікацій.

Інформація про авторів:

Мокрицька Ольга Володимирівна, канд. техн. наук, доцент, кафедра систем автоматизованого проектування.

Email: olha.v.mokrytska@lpnu.ua; <https://orcid.org/0000-0002-2887-9585>

Мочернюк Юрій Миколайович, аспірант, кафедра інженерії програмного забезпечення.

Email: mocherniuk.iurii@nltu.lviv.ua; <https://orcid.org/0009-0002-1370-363X>

Цитування за ДСТУ: Мокрицька О. В., Мочернюк Ю. М. Використання алгоритмів машинного навчання для автоматизації процесу модераторів контенту в групових чатах месенджерів. Науковий вісник НЛТУ України. 2024, т. 34, № 7. С. 52–59.

Citation APA: Mokrytska, O. V., & Mocherniuk, Yu. M. (2024). Application of machine learning algorithms for automation of content moderation in app group messaging chats. *Scientific Bulletin of UNFU*, 34(7), 52–59. <https://doi.org/10.36930/40340707>

Об'єкт дослідження – використання алгоритмів машинного навчання.

Предмет дослідження – методи і засоби машинного навчання, які дають змогу автоматизувати процес модерації контенту в групових чатах месенджерів.

Мета роботи – дослідити особливості використання відомих алгоритмів машинного навчання для автоматизації процесу модерації контенту в групових чатах месенджерів, що дасть змогу інтегрувати методи вибору функцій класифікації з класифікаційними моделями.

Для досягнення зазначеної мети визначено такі основні завдання дослідження:

1. Проаналізувати наявні системи модерації на підставі машинного навчання, що дасть змогу виявити серед них придатні алгоритми та покращити ефективність автоматизації процесу модерації контенту.
2. Інтегрувати алгоритми вибору функцій класифікації, такі як Particle Swarm Optimization (PSO), Salp Swarm Algorithm (SSA) і Grey Wolf Optimization (GWO), в автоматизованому процесі модерації, що дасть можливість підвищити точність класифікації контенту, зменшити обчислювальні витрати та покращити ефективність виявлення небажаних або шкідливих матеріалів у реальному часі.
3. Оцінити ефективність моделей класифікації, зокрема Support Vector Machines (SVM) і k-Nearest Neighbors (kNN), щодо виявлення неприйнятної контенту, які допоможуть визначити найкращі підходи для підвищення точності та швидкості модерації, мінімізувати кількість помилкових спрацювань та оптимізувати автоматизовані системи для ефективного оброблення великих обсягів даних.
4. Провести порівняльний аналіз результативності обраних алгоритмів за ключовими показниками (точність, чутливість, специфічність тощо), що забезпечить обґрунтований вибір найефективнішої моделі для автоматизації процесу модерації контенту, підвищить надійність виявлення неприйнятної матеріалу та зменшить кількість хибнопозитивних і хибнонегативних результатів.

Виконання цих завдань дасть змогу створити ефективну, адаптивну систему модерації, яка підвищить точність виявлення шкідливого контенту та зменшить кількість помилкових спрацювань, що значно покращить якість комунікацій в онлайн-чатах.

Аналіз останніх досліджень та публікацій. Останні наукові дослідження демонструють важливість впровадження алгоритмів машинного навчання для автоматизації процесу модерації контенту на цифрових платформах, зокрема у групових чатах месенджерів. У роботі [9] вказано, що через значний обсяг створюваного користувачами контенту ручна модерація стає практично неможливою. Сучасні підходи, такі як використання штучного інтелекту, зокрема – методів оброблення природної мови, а саме – нейролінгвістичного програмування NLP (англ. *Natural Language Processing*) для аналізу тексту та комп'ютерного зору для модерації візуальних даних, є основними інструментами для автоматизації цього процесу. Це дає змогу виявляти шкідливий або неприйнятний контент у реальному часі, значно зменшуючи навантаження на модераторів і підвищуючи точність відсіву. Важливо зазначити, що групові чати мають динамічний характер, де швидкі зміни у тематиці та стилі спілкування потребують адаптивних алгоритмів, здатних оперативного реагувати на нові виклики.

В іншому дослідженні, присвяченому питанням виявлення кіберзлочинів за допомогою класифікаційної

моделі на підставі методу опорних векторів SVM (англ. *Support Vector Machine*) і моделі на підставі методу k-найближчих сусідів kNN (англ. *k-Nearest Neighbor Method*) [12], проаналізовано ефективність таких алгоритмів для виявлення потенційно небезпечного контенту, зокрема щодо злочинної діяльності в онлайн-комунікаціях. Ці підходи також можна застосувати в системах автоматизованого процесу модерації контенту для відстеження кібербулінгу та інших форм онлайн-агресії в месенджерах.

Дослідження [14] звертає увагу на створення унікального набору даних NoisyHate, який містить спеціально зібрані та очищені збурення, написані людьми, з різних онлайн-платформ. Зібрані збурення було використано для тестування сучасних мовних моделей, таких як BERT і RoBERTa, які показали, що навіть передові алгоритми машинного навчання мають труднощі з обробленням таких атипових даних. Ці результати вказують на потребу розроблення нових, більш стійких, моделей для модерації контенту в умовах реальних людських помилок, та неточностей. Зокрема, моделі, такі як BERT і RoBERTa, мають потенціал для вдосконалення саме щодо групових чатів, де часті помилки або неправильні висловлювання можуть швидко призвести до загострення конфліктів.

Автори дослідження [3] у своїй роботі розглянули поточні досягнення у сфері модерації соціальних медіа. Автори зосереджують увагу на важливості використання систем на підставі штучного інтелекту для виявлення шкідливого контенту, в тому числі образливих повідомлень, спаму і фейкових новин. Однією з основних проблем, з якою стикаються сучасні системи модерації, є складність адаптації алгоритмів до змінюваного контенту і виявлення контенту, що порушує етичні або юридичні норми. Також потрібно зазначити, що етичні виклики, пов'язані з модерацією контенту в групових чатах, часто мають більш складний характер через приватність комунікацій та ризик потенційних порушень конфіденційності.

У роботі [13] розглянуто підхід до інкрементального машинного навчання для класифікації тексту в системах модерації коментарів. Ця технологія дає змогу системам покращувати свою ефективність з часом, адаптуючись до нових викликів, що виникають у процесі аналізу контенту в режимі реального часу.

Окрім цього, дослідження [1, 2] зосереджуються на питаннях безпеки мереж і виявлення атак, таких як DDoS (англ. *Distributed Denial-Of-Service* – розподілена відмова в обслуговуванні) та DoS (англ. *Denial-Of-Service* – відмова в обслуговуванні) з використанням алгоритмів машинного навчання. Незважаючи на те, що ці дослідження мають більш технічний характер, запропоновані в них підходи до вибору ознак і оптимізації моделей можна адаптувати для модерації контенту в месенджерах, зокрема для виявлення шкідливої активності або спаму.

Дослідження [4] детально розглядає застосування різноманітних технік машинного навчання, в т.ч. класифікацію та кластеризацію, для вдосконалення процесів автоматичної модерації. Важливо, що автори вказують на потребу постійного вдосконалення алгоритмів для підвищення точності виявлення неприпустимого контенту в реальному часі та зниження кількості хибнопозитивних результатів.

Загалом аналіз сучасних публікацій демонструє важливість використання складних і адаптивних алгоритмів машинного навчання для ефективної автоматизації процесу модерації контенту. Оскільки середовище онлайн-спілкування постійно змінюється, важливо забезпечувати постійний моніторинг та оновлення алгоритмів, щоб вони залишалися ефективними і відповідали новим викликам, що виникають у процесі модерації контенту в групових чатах месенджерів. Попри значний прогрес, залишається низка викликів, зокрема, щодо забезпечення етичності, справедливості та прозорості заснованих систем.

Матеріали та методи дослідження. У дослідженні використано методи машинного навчання, зокрема алгоритми вибору функцій класифікації (Particle Swarm Optimization, Salp Swarm Algorithm, Grey Wolf Optimization) та класифікаційні моделі (Support Vector Machines, k-Nearest Neighbors), для автоматизації процесу модерації контенту в групових чатах Telegram. Для аналізу ефективності застосовано стандартні метрики продуктивності: точність, чутливість, специфічність, вивірненість і оцінка F1, опис обчислення яких наведено далі. Для моделювання, як джерело даних, було взято підмножину з відкритих наборів даних "Hate Speech and Offensive Language" та "Jigsaw Unintended Bias in Toxicity Classification" обсягом в 5000 прикладів для об'єктивної оцінки алгоритмів. Розрахунки виконано шляхом усереднення значень кожної метрики по всіх наборах даних, із застосуванням крос-валідації для підвищення достовірності результатів. Точність отриманих значень підтверджена обчисленням довірчих інтервалів, що демонструють стабільність метрик на різних вибірках. Усі алгоритми та моделі реалізовані за допомогою Python-бібліотек Scikit-learn та PySwarm.

Результати дослідження та їх обговорення / Research results and their discussion

Дослідження ефективності алгоритмів машинного навчання проводиться через Telegram-бота, що модерує групи, видаляє неприйнятні повідомлення, блокує користувачів і здійснює аналіз настроїв. Бот написаний на Python з використанням кількох бібліотек і PostgreSQL для реєстрації всіх дій. Налаштування починається зі створення віртуального середовища для ізоляції залежностей та збереження URL-адреси бази даних у змінній TELEGRAM_BOT_POSTGRES_URL.

Змінні середовища керують функціональністю чатів Telegram: MESSAGE_BAN_PATTERNS і MESSAGE_HIDE_PATTERNS містять регулярні вирази для блокування чи видалення повідомлень; NAME_BAN_PATTERNS – для блокування користувачів за іменем. Бот може моніторити кілька чатів через змінну CHAT_IDS і має режим налагодження (DEBUG), де дії відображаються безпосередньо в чаті.

Функція NOTIFY_CHAT надсилає звіти про дії бота в окремий чат. Бот також підтримує команду / price через CMC_API_KEY для отримання цін на криптовалюти. Налаштування шаблонів для виявлення небажаних повідомлень дає змогу точну модерацію. Бот може працювати як локально, так і на платформах на кшталт Heroku, де його функціонування активується через інформаційну панель.

Автоматизація процесу модерації містить використання алгоритмів вибору функцій класифікації, таких

як Gray Wolf Optimization (GWO), Particle Swarm Optimization (PSO) і Salp Swarm Algorithm (SSA) у поєднанні з такими моделями класифікації, як Support Vector Machine (SVM) і k-Nearest Neighbors (kNN). Ці алгоритми допомагають боту виявляти значні шаблони в повідомленнях чату, оптимізуючи процес модерації та точно класифікуючи вміст.

Набори даних, які було використано для основного навчання бота, містять журнали повідомлень із кількох групових чатів Telegram, які попередньо позначені як "звичайні" або "позначені" на підставі попередньо визначених критеріїв модерації. Кожен запис містить такі функції, як:

- вміст повідомлення (текст);
- ідентифікатор користувача та ім'я користувача;
- мітки часу;
- довжина повідомлення;
- частота певних ключових слів або шаблонів.

Названі функції екстрагують та зберігаються в базі даних. Для навчання загальнодоступні набори даних, що містять дані виявлення спаму, такі як "Набір даних про спам для SMS" і "Набір даних про спам Enron", можна адаптувати та об'єднати з даними Telegram для створення повного набору даних.

На фрагменті коду 1 показано, як екстрагують функції з повідомлень, що зберігаються в базі даних, щоб підготувати їх до оброблення. Функція extract_features_from_db отримує всі повідомлення з бази даних за допомогою сеансу SQLAlchemy. Для кожного повідомлення обчислюються основні параметри:

- msg_len: довжина повідомлення.
- uppercase_count: кількість літер у верхньому регістрі.
- contains_link: двійковий показник (1 або 0), якщо повідомлення містить URL-адресу.

Повідомлення позначаються як "звичайні" (0) або "позначені" (1) на підставі наявних даних. Цей екстрагований набір даних, вказаний як x (об'єкти) та y (мітки), використовуватиметься для навчання вибору об'єктів і алгоритмів класифікації.

Вибір класифікаційних функцій має вирішальне значення для зменшення шуму та підвищення ефективності процесу класифікації бота. Якщо така функція обмежена на множині і зверху і знизу, то її називають обмеженою на всій множині наявних даних. Фрагменти коду 2-4 ілюструють, як алгоритми PSO, SSA та GWO застосовують для вибору найбільш релевантних функцій із даних повідомлень.

У фрагменті коду 2 реалізовано алгоритм PSO для вибору ознак. У стадії ініціалізації 30 частинок генерують випадково, де кожна частинка вказує на потенційну підмножину ознак (двійковий масив, де 1 вказує на вибрану функцію). При ітераційному циклі в понад 100 ітерацій піднабір функцій кожної частинки оцінюється за допомогою логіки класу SelectKBest для обчислення оцінки придатності. Далі частинки коригують свої позиції в просторі функцій на підставі власного досвіду (особистого рекорду) і випадкових впливів, імітуючи, як частинки "рояться" до оптимальних рішень. Внаслідок повертається найефективніша підмножина функцій. Алгоритм PSO допомагає зменшити розмірність даних, даючи можливість боту зосереджуватися тільки на найбільш інформативних особливостях повідомлень, цим самим покращуючи ефективність і точність подальшого процесу класифікації.

```

import numpy as np
from sklearn.feature_selection import SelectKBest, f_classif
def extract_features_from_db():
    session = session()
    messages = session.query(Message).all()
    features = []
    labels = []
    for message in messages:
        msg_len = len(message.message)
        uppercase_count = sum(1 for c in message.message if c.isupper())
        contains_link = 1 if re.search(r"http[s]?://", message.message) else 0
        features.append([msg_len, uppercase_count, contains_link])
        labels.append(1 if "flagged" in message.message.lower() else 0)
    session.close()
    return np.array(features), np.array(labels)
X, y = extract_features_from_db()

```

Код 1. Екстракція функцій з бази даних / Extraction of features from the database

```

def pso_feature_selection(X, y):
    num_particles = 30
    num_features = X.shape[1]
    particles = np.random.randint(2, size=(num_particles, num_features))
    velocities = np.zeros((num_particles, num_features))
    personal_best_positions = particles.copy()
    personal_best_scores = np.zeros(num_particles)
    for iteration in range(100):
        for i in range(num_particles):
            selected_features = particles[i]
            selected_X = X[:, selected_features == 1]
            if selected_X.shape[1] == 0:
                continue
            k_best = SelectKBest(score_func=f_classif, k='all').fit(selected_X, y)
            score = k_best.scores_.sum()
            if score > personal_best_scores[i]:
                personal_best_scores[i] = score
                personal_best_positions[i] = particles[i]
                velocities[i] += 0.5 * (personal_best_positions[i] - particles[i]) + 0.3 * np.random.rand(num_features)
                particles[i] = np.where(np.random.rand(num_features) < velocities[i], 1, 0)
        best_particle_index = np.argmax(personal_best_scores)
        return np.where(personal_best_positions[best_particle_index] == 1)[0]
selected_features_pso = pso_feature_selection(X, y)
print("Selected features using PSO:", selected_features_pso)

```

Код 2. Вибір класифікаційних функцій за допомогою алгоритму PSO / Selection of classification features using PSO algorithm

Фрагмент коду 3 демонструє роботу алгоритму SSA у задачі вибору функцій класифікації сальпів. При цьому створюється рій сальпів (20 агентів), кожен з яких вказує на потенційну підмножину відповідних функцій. Понад 50 ітерацій сальпи слідує за рухами лідера з

коригуванням, внесеним на підставі параметра конвергенції (c1). Це імітує те, як сальпи досліджують простір функцій. Найоптимальніша сальпа на підставі фітнес-оцінки стає лідером, спрямовуючи інших сальп до кращих комбінацій функцій.

```

def ssa_feature_selection(X, y):
    num_salps = 20
    num_features = X.shape[1]
    salps = np.random.randint(2, size=(num_salps, num_features))
    leader = salps[0]
    for iteration in range(50):
        for i in range(1, num_salps):
            for j in range(num_features):
                c1 = 2 * np.exp(-(4 * iteration / 50)**2)
                if j == 0:
                    leader[j] = np.random.randint(0, 2)
                salps[i, j] = (leader[j] + c1 * (np.random.rand() - 0.5)) % 2
            selected_features = np.where(salps[i] == 1)[0]
            if len(selected_features) > 0:
                score = SelectKBest(score_func=f_classif, k='all').fit(X[:, selected_features], y).scores_.sum()
                current_score = SelectKBest(score_func=f_classif, k='all').fit(X[:, leader == 1], y).scores_.sum()
                if score > current_score:
                    leader = salps[i].copy()
        return np.where(leader == 1)[0]
selected_features_ssa = ssa_feature_selection(X, y)
print("Selected features using SSA:", selected_features_ssa)

```

Код 3. Вибір класифікаційних функцій за допомогою алгоритму SSA / Selection of classification functions using SSA algorithm

Алгоритм GWO, показаний у фрагменті коду 4, оптимізує процес вибору функцій класифікації. Ієрархія вовків: вовки (альфа, бета, дельта) представляють різні потенційні рішення, причому альфа є найоптимальнішим. Упродовж понад 50 ітерацій вовки коригують свої

```
def gwo_feature_selection(X, y):
    alpha, beta, delta = np.random.randint(2, size=X.shape[1]),
        np.random.randint(2, size=X.shape[1]),
        np.random.randint(2, size=X.shape[1])
    wolves = [alpha, beta, delta]
    for iteration in range(50):
        a = 2 - iteration * (2 / 50)
        for wolf in wolves:
            for i in range(X.shape[1]):
                d_alpha = abs(2 * np.random.rand() * alpha[i] - wolf[i])
                d_beta = abs(2 * np.random.rand() * beta[i] - wolf[i])
                d_delta = abs(2 * np.random.rand() * delta[i] - wolf[i])
                wolf[i] = (d_alpha + d_beta + d_delta) / 3
    return np.where(alpha == 1)[0]
selected_features_gwo = gwo_feature_selection(X, y)
print("Selected features using GWO:", selected_features_gwo)
```

Код 4. Вибір класифікаційних функцій за допомогою алгоритму GWO / Selection of classification features using GWO algorithm

Рішення використання ансамблевої моделі у проти- вагу індивідуальному застосуванню алгоритмів вибору функцій класифікації зумовлено тим, що алгоритм PSO швидко об'єднується для виявлення очевидних особливостей, алгоритм SSA збалансовує розвідку та експлуатацію для уточнення вибору, а алгоритм GWO перевершує пошук складних багатовимірних моделей. Алгоритм PSO моделює соціальну поведінку частинок і є ефективним у швидкому переході до оптимального рішення. Він особливо сильний для ефективного дослідження простору пошуку, але може потрапити в пастку локальних оптимумів, коли складність проблеми зростає. Алгоритм SSA, натхненний харчовою поведінкою сальп в океані, пропонує баланс між розвідкою (пошук різноманітних рішень) і експлуатацією (удосконалення наявних рішень). Це дає змогу ефективно уникнути передчасної конвергенції, але іноді може призводити до

```
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
X_selected = X[:, selected_features_pso]
X_train, X_test, y_train, y_test = train_test_split(X_selected, y, test_size=0.3, random_state=42)
svm_model = SVC(kernel='linear', C=1.0)
svm_model.fit(X_train, y_train)
print(f"SVM Accuracy: {svm_model.score(X_test, y_test)}")
```

Код 5. Класифікація повідомлень із використанням моделі SVM / Message classification by utilizing SVM model

Класифікатор KNeighborsClassifier навчається за допомогою того самого навчального набору даних. Його точність оцінюється на тестовому наборі.

```
from sklearn.neighbors import KNeighborsClassifier
knn_model = KNeighborsClassifier(n_neighbors=3)
knn_model.fit(X_train, y_train)
print(f"KNN Accuracy: {knn_model.score(X_test, y_test)}")
```

Код 6. Класифікація повідомлень із використанням моделі kNN / Message classification by utilizing the kNN model

У фрагменті коду 7 наведено модель повністю інтегрованого та скомпільованого алгоритму, асистованого методами машинного навчання ML для модерації повідомлень Telegram-ботом. Клас EnhancedTelegramMonitorBot розширює наявного бота, в т.ч. навчені svm_model і knn_model. Метод moderate_message перевіряє вхідні повідомлення на прогнози поведінки моделей і вживає відповідних заходів, наприклад, видаляючи позначені повідомлення. Спеціальний метод extract_featu-

позиції на підставі позицій альфа-, бета- та дельта-вовків, поступово наближаючись до оптимальної підмножини функцій. Внаслідок роботи алгоритму остаточна позиція альфа-вовка вказує на найефективнішу підмножину функцій для класифікації повідомлень.

деякого сповільнення пошуку найкращого рішення виходу з ситуації, що склалася. Алгоритм GWO імітує лідерство та мисливську поведінку сірих вовків і врівноважує дослідження та експлуатацію.

Наступним необхідним кроком є класифікація повідомлень залежно від їх визначених властивостей. Ця операція виконується з використанням класифікаційних моделей SVM та kNN, використання яких проілюстроване у фрагментах коду 5-6.

У фрагменті коду 5 показано процес поділу набору даних на набори для навчання та тестування з використанням функцій, вибраних алгоритмом PSO. Модель SVC навчається за допомогою лінійного ядра, підлаштовуючи модель для того, щоб відрізнити звичайні повідомлення від позначених. Визначається точність моделі, що вказує на її ефективність у передбаченні того, чи потрібно позначати повідомлення.

ges забезпечує використання правильних функцій для класифікації даних.

Для досягнення мети дослідження використовують метрики продуктивності, які надають кількісні засоби для вимірювання здатності бота правильно ідентифікувати неприйнятні повідомлення та керувати ними на підставі його можливостей класифікації. Перший показник – точність, вимірює, наскільки прогнози бота узгоджують з фактичними результатами модерації, розраховується як відношення правильно класифікованих повідомлень (як позначених, так і не позначених) до всіх оброблених повідомлень. Однак, через можливий дисбаланс використаних класів, де неприйнятні повідомлення можуть надходити рідше, ніж звичайні повідомлення, сама по собі точність може не повністю відобразити ефективність бота.

Чутливість модерації вмісту вказує на частку правильно визначених неприйнятних повідомлень серед

усіх фактично позначених повідомлень, що має ви-
 шальне значення для того, щоб бот зафіксував най-
 більш проблематичний вміст. Специфічність модерації
 вмісту вимірює частку звичайних повідомлень, точно
 визначених як невідповідні, допомагаючи цим самим
 мінімізувати помилкові спрацьовування та гарантувати,

```
class EnhancedTelegramMonitorBot(TelegramMonitorBot):
    def __init__(self):
        super().__init__()
        self.svm_model = svm_model
        self.knn_model = knn_model
    def moderate_message(self, bot, update):
        message_features = self.extract_features(update.message.text)
        prediction_svm = self.svm_model.predict([message_features])
        prediction_knn = self.knn_model.predict([message_features])
        if prediction_svm == 1 or prediction_knn == 1:
            update.message.delete()
            log_violation(update.message)
    def extract_features(self, message_text):
        return [len(message_text), sum(1 for c in message_text if c.isupper()),
                int(bool(re.search(r"http[s]?://", message_text)))]
c = EnhancedTelegramMonitorBot()
c.start()
```

Код 7. Інтегрований та зкомпільований алгоритм асистованої модерації повідомлень / Integrated and compiled algorithm of assisted message moderation

Оцінка F1 є гармонійним середнім значенням точ-
 ності та запам'ятовування, забезпечує збалансоване уяв-
 лення про продуктивність бота, особливо коли потрібно
 мінімізувати помилкові спрацьовування та помилкові
 негативні результати одночасно. Рисунок демонструє
 порівняння різних комбінацій алгоритмів залежно від
 метрик продуктивності.

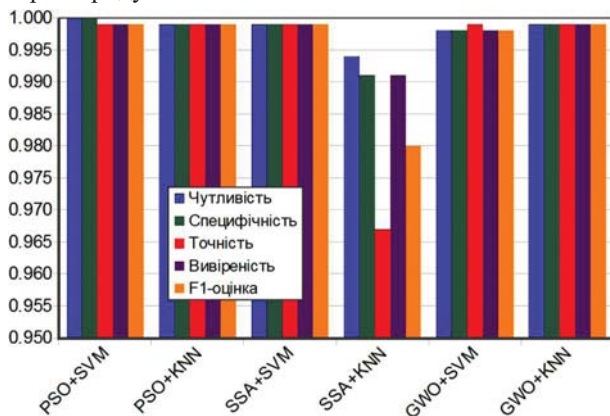


Рисунок. Порівняння метрик продуктивності / Comparison of performance metrics

Показники чутливості та специфічності демонстру-
 ють незмінно високі значення в усіх моделях, що вказує
 на те, що кожна модель однаково ефективна у визна-
 ченні справжніх позитивних і справжніх негативних ре-
 зультатів, зберігаючи свою ефективність близько до ма-
 ксимального значення одиниця. Однак, спостерігається
 значне відхилення в показнику точності для моделі
 SSA+kNN, яка відчуває різке зниження порівняно з ін-
 шими. Це падіння означає, що модель SSA+kNN більше
 намагається відрізнити справжні позитивні результати
 від загальної кількості прогнозованих неприйнятних по-
 відомлень, що робить його менш надійним у цьому зна-
 ченні. Решта моделей зберігають високу точність, де-
 монструючи високу ефективність. Точність і оцінка
 F1 залишаються відносно стабільними для всіх моде-
 лей, підтверджуючи їх загальну збалансовану продук-
 тивність у класифікаційних завданнях. Отже, хоча біль-
 шість моделей демонструють послідовну та високу

що законні повідомлення не позначаються несправед-
 ливо. Вивіренисть модерації вмісту оцінює, скільки по-
 відомлень, позначених ботом, є насправді неприйнятни-
 ми. Цей показник дає уявлення про ефективність робо-
 ти бота щодо уникнення непотрібних дій проти закон-
 них користувачів.

ефективність, явне зниження точності моделі SSA+kNN
 є помітним винятком, який підкреслює його потенційну
 слабкість у точності прогнозування.

Додатковий аналіз, результати якого наведено у таб-
 лиці, зосереджував увагу на кількості функцій, вибра-
 них кожним алгоритмом на етапі вибору як таких.

Таблиця. Вибрані функції за різними алгоритмами /
 Selected features based on various algorithms

Алгоритм	PSO+ SVM	PSO+ kNN	SSA+ SVM	SSA+ kNN	GWO+ SVM
Загальна кількість роз- глянутих ознак	80	80	80	80	80
Відібрані ознаки	35	38	28	25	42
Ефективність вибору оз- нак, %	43.75	47.50	35.00	31.25	52.50
Видалені зайві ознаки	45	42	52	55	38
Тривалість виконання, мс	150	170	140	130	160

Комбінацією алгоритму PSO+SVM було вибрано 35
 ознак, тоді як алгоритмом GWO+kNN – 42 ознаки, що
 демонструє концепцію використання різних алгорит-
 мів, які віддають пріоритет різним особливостям вмісту
 повідомлення під час визначення шаблонів, що вказу-
 ють на невідповідну поведінку. Ці варіації свідчать про
 те, що, хоча деякі комбінації алгоритмів є більш вибір-
 ковими, інші є комплексними, забезпечуючи ширше
 охоплення потенційно шкідливих ознак.

Обговорення результатів дослідження. Викорис-
 тання алгоритмів вибору функцій класифікації у проце-
 сі автоматизованого процесу модерації повідомлень бу-
 ло ретельно досліджено в різних роботах. Наприклад, у
 дослідженні [4] розглянуто застосування машинного
 навчання для модерації соціальних медіа, зокрема для
 виявлення шкідливого контенту. Їхнє дослідження по-
 казало, що традиційні методи, засновані на фіксованих
 правилах, істотно поступаються алгоритмам машинно-
 го навчання у значенні точності, що відповідає резуль-
 татам нашого дослідження, де комбінація алгоритмів
 PSO+SVM досягла найвищої точності.

В іншій роботі [6] також досліджували проблему ви-
 бору функцій класифікації щодо аналізу великих даних

і використання алгоритмів оптимізації. У своєму дослідженні автор дійшов висновку, що алгоритми оптимізації, такі як алгоритм GWO, можуть значно підвищити ефективність вибору функцій класифікації. У нашому випадку використання алгоритму GWO показало позитивні результати, однак алгоритм PSO виявився більш стабільним і ефективним у середовищі Telegram.

У роботі [5] автори розглянули дещо інші методи класифікації даних, такі методи, як Decision Trees і Random Forests, призначені для виявлення шкідливого контенту в онлайн-середовищах. Хоча ці методи продемонстрували прийнятну продуктивність, у нашому ж дослідженні алгоритм SVM виявився більш точним порівняно з алгоритмом Random Forests, особливо в поєднанні з алгоритмом PSO.

Останні дослідження в галузі автоматизованого процесу модерації, наприклад робота [10], підтверджують важливість вибору функцій класифікації для зменшення розмірності даних і підвищення точності класифікації. Їхній підхід до використання алгоритмів вибору функцій класифікації, таких як модифікації PSO, виявився ефективним, що корелює з нашими висновками, де алгоритм PSO показав високу здатність адаптуватися до змін у повідомленнях.

Окрім цього, автори у дослідженні [11] розглянули поєднання моделей класифікації, таких як kNN, з різними методами оптимізації. Їхні результати показали, що kNN у поєднанні з алгоритмами вибору функцій класифікації може покращити продуктивність, проте ми виявили, що комбінація алгоритмів SSA+kNN мала меншу стабільність порівняно з алгоритмом PSO+SVM, особливо у різних умовах.

Загалом наші результати підтверджують висновки інших дослідників про ефективність застосування алгоритмів вибору функцій класифікації у поєднанні з класифікаційними моделями, однак важливою перевагою нашого підходу є збалансованість між точністю та швидкістю, яку забезпечує алгоритм PSO.

Отже, внаслідок виконаної роботи можна сформулювати такі наукову новизну та практичну значущість результатів дослідження.

Наукова новизна отриманих результатів дослідження – розроблено метод автоматизованого процесу модерації контенту у групових чатах месенджерів, який інтегрує алгоритми вибору функцій класифікації (Particle Swarm Optimization, Salp Swarm Algorithm та Grey Wolf Optimization) з класифікаційними моделями Support Vector Machines і k-Nearest Neighbors, що дає змогу значно підвищити точність виявлення шкідливого контенту, зменшити кількість помилкових спрацьовувань та забезпечити адаптивність системи до мінливого середовища комунікацій.

Практична значущість результатів дослідження – розроблено та впроваджено автоматизовану систему модерації у групових чатах месенджерів, зокрема для підвищення ефективності фільтрації шкідливого контенту шляхом використання алгоритмів вибору функцій класифікації (Particle Swarm Optimization, Salp Swarm Algorithm, Grey Wolf Optimization) та класифікаційних моделей (Support Vector Machines і k-Nearest Neighbors), що дає змогу підвищити точність і адаптивність системи, забезпечити дещо безпечніше комунікаційне середовище, а також потенційно масштабувати рішення для інших платформ та доменів.

Висновки / Conclusions

Досліджено особливості використання відомих алгоритмів машинного навчання для автоматизації процесу модерації контенту в групових чатах месенджерів, що дало змогу інтегрувати методи вибору функцій класифікації з класифікаційними моделями, підвищити ефективність і точність ідентифікації та усунути небажаний вміст. За результатами проведеного дослідження можна зробити такі основні висновки.

1. Систематизовано підходи до автоматизованого процесу модерації контенту на підставі машинного навчання, що поєднують алгоритми вибору функцій класифікації (PSO, SSA, GWO) з класифікаційними моделями (SVM, kNN), що дає змогу підвищити ефективність виявлення шкідливого контенту.
2. Проаналізовано продуктивність різних алгоритмів вибору функцій класифікації щодо їх використання для модерації чатів. Встановлено, що алгоритм PSO забезпечує найбільш стабільні результати модерації чатів завдяки швидкій адаптації до мінливого комунікаційного середовища.
3. Встановлено, що комбінація алгоритмів PSO+SVM досягла найвищої точності у класифікації повідомлень, перевершуючи інші алгоритми завдяки збалансованості між швидкістю та точністю процесу модерації.
4. Досліджено варіанти використання алгоритмів SSA та GWO у поєднанні з моделлю класифікації kNN для модерації чатів, що показали ефективність у різних умовах, але продемонстрували меншу стабільність порівняно з алгоритмом PSO.
5. Проведено оцінювання результативності роботи алгоритмів за допомогою таких показників, як точність, чутливість, специфічність, вивіреність та оцінка F1, що дало можливість підтвердити значний вплив вибору функцій класифікації на точність модерації.
6. Наведено потенціал методів і алгоритмів машинного навчання для подальшого масштабування та застосування запропонованої системи модерації на інших платформах, що може сприяти підвищенню безпеки обміну інформацією в онлайн-комунікаціях загалом.

References

1. Ahmad, R., Wazirali, R., Bsoul, Q., Abu-Ain, T., & Abu-Ain, W. (2021). Feature-selection and mutual-clustering approaches to improve DoS detection and maintain WSNs lifetime. *Sensors*, 21(14). <https://doi.org/10.3390/s21144821>
2. Almaiah, M., Almaiah, D., Alrawashdeh, R., Alkhdour, T., Al-Ali, R., Rjoub, G., & Aldahyani, T. (2024). Detecting DDoS attacks using machine learning algorithms and feature selection methods. *International Journal of Data and Network Science*, 8(2), 45–58. <https://doi.org/10.5267/j.ijdns.2024.6.001>
3. Gongane, V., Munot, M., & Anuse, D. (2022). Detection and moderation of detrimental content on social media platforms: current status and future directions. *Social Network Analysis and Mining*, 12. <https://doi.org/10.1007/s13278-022-00951-3>
4. Gulati, G., Jha, H. A., Jain, R., Sharma, M., & Chaudhary, V. (2024). Content moderation system using machine learning techniques. In: Hassanien, A. E., Castillo, O., Anand, S., & Jaiswal, A. (Eds.). *International Conference on Innovative Computing and Communications*. Springer, Singapore. https://doi.org/10.1007/978-981-99-4071-4_58
5. Han, H., Asif, M., Awwad, E. M., et al. (2024). Innovative deep learning techniques for monitoring aggressive behavior in social media posts. *Journal of Cloud Computing*, 13, 19 p. <https://doi.org/10.1186/s13677-023-00577-6>
6. Hongyu, P., Shanxiong, C., & Hailing, X. (2023). A high-dimensional feature selection method based on modified Gray Wolf Optimization. *Applied Soft Computing*, 135. <https://doi.org/10.1016/j.asoc.2023.110031>

7. Hrytsiuk, Y. I. (2022). Features of giving preference to the characteristics of the software product quality model. *Scientific Bulletin of UNFU*, 32(3), 79–102. <https://doi.org/10.36930/40320313>
8. Liu, X., & Du, Y. (2023). Towards effective feature selection for IoT botnet attack detection using a genetic algorithm. *Electronics*, 12(5). <https://doi.org/10.3390/Electronics12051260>
9. Mantri, A. (2021). Real-Time Content Moderation Using Artificial Intelligence and Machine Learning. *International Journal of Scientific and Engineering Research*, 10, 1682–1684. <https://doi.org/10.21275/SR24724150350>
10. Nachaoui, M., Lakouam, I., & Hafidi, I. (2024). Hybrid particle swarm optimization algorithm for text feature selection problems. *Neural Comput & Applic*, 36, 7471–7489. <https://doi.org/10.1007/s00521-024-09472-w>
11. Sang, B., Xu, W., Chen, H., & Li, T. (2023). Active antinoise fuzzy dominance rough feature selection using adaptive k-Nearest Neighbors. *IEEE Transactions on Fuzzy Systems*, 31(11), 3944–3958. <https://doi.org/10.1109/TFUZZ.2023.3272316>
12. Veena, K., Meena, K., Teekaraman, Y., Kuppusamy, R., & Radhakrishnan, A. (2022). SVM Classification and kNN Techniques for Cyber Crime Detection. *Security Threats and Challenges in Future Mobile Communication Systems*. <https://doi.org/10.1155/2022/3640017>
13. Wolters, A., Müller, K., & Riehle, D. M. (2022). Incremental Machine Learning for Text Classification in Comment Moderation Systems. *Lecture Notes in Computer Science*, 13427, 160–171. https://doi.org/10.1007/978-3-031-18253-2_10
14. Ye, Y., Le, T., & Lee, D. (2023). NoisyHate: Benchmarking Content Moderation Machine Learning Models with Human-Written Perturbations Online. *Computer Science > Machine Learning*. arXiv:2303.10430v1 <https://doi.org/10.48550/arXiv.2303.10430>

O. V. Mokrytska¹, Yu. M. Mocherniuk²

¹ Lviv Polytechnic National University, Lviv, Ukraine

² Ukrainian National Forestry University, Lviv, Ukraine

APPLICATION OF MACHINE LEARNING ALGORITHMS FOR AUTOMATION OF CONTENT MODERATION IN APP GROUP MESSAGING CHATS

The integration of machine learning (ML) methods into an automated moderation system for Telegram group chats aimed at improving large-scale communication management has been analyzed. In the course of research, the moderation of large group chats is identified to become especially challenging due to the increasing volume of messages and content diversity, necessitating an efficient filtering system to ensure a safe communication environment. Advanced feature selection algorithms, such as Particle Swarm Optimization (PSO), Salp Swarm Algorithm (SSA), and Gray Wolf Optimization (GWO), have been applied to automate the detection of the most relevant message characteristics, significantly enhancing the moderation process. The use of classification models, including Support Vector Machines (SVM) and k-Nearest Neighbors (kNN), for identifying inappropriate or harmful content has been studied. It has also been established that feature selection is critical for reducing data dimensionality and improving classification accuracy. A comparative analysis of the effectiveness of PSO, SSA, and GWO in the context of chat moderation has been conducted. It has been shown that PSO demonstrates high efficiency due to its rapid adaptation to the changing communication environment, while SSA and GWO exhibit variations in their approaches to global optimization of feature selection. A classification system has been developed in which the PSO+SVM combination achieves the highest accuracy, outperforming other algorithms in balancing speed and accuracy. Results for SSA+kNN and GWO+kNN combinations have also been presented, suggesting high effectiveness but less stable performance in various conditions. The performance of the algorithms has been evaluated using metrics such as accuracy, sensitivity, specificity, precision, and F1 score. The comprehensive analysis confirmed that the implementation of advanced feature selection algorithms significantly improves the accuracy of harmful content detection, reducing false positives and enhancing the overall efficiency of the moderation process. The study of the respective algorithms in the context of moderation also identified the potential for deploying the developed system on other platforms and expanding into new domains to provide a scalable and adaptable solution for maintaining safety in online communications.

Keywords: Telegram bot; PSO; SSA; GWO; SVM; feature selection; classification; ensemble methods.