

Батечко Д. О.,
здобувач вищої освіти спеціальності 122 Комп'ютерні науки
Науковий керівник: Богатенкова О.Є., асистент кафедри економічної
кібернетики, комп'ютерних наук та інформаційних технологій
Миколаївський національний аграрний університет
м. Миколаїв

ПОРІВНЯЛЬНИЙ АНАЛІЗ МЕТОДІВ КЛАСИФІКАЦІЇ ДЛЯ ВИЯВЛЕННЯ СПАМУ В ІНФОРМАЦІЙНИХ СИСТЕМАХ ОРГАНІВ ДЕРЖАВНОЇ ВЛАДИ

Стрімке зростання обсягів текстової інформації в цифрових комунікаційних середовищах зумовлює актуальність задачі автоматичного виявлення та фільтрації спаму. Небажані повідомлення не лише знижують якість користувацького досвіду, а й можуть становити загрозу інформаційній безпеці, зокрема через фішингові атаки та поширення шкідливого програмного забезпечення. У зв'язку з цим застосування методів машинного навчання для класифікації текстових даних набуває особливої практичної значущості.

Особливої актуальності проблема спам-фільтрації набуває в інформаційних системах органів державної влади, де електронна пошта, онлайн-звернення громадян, системи електронного документообігу та державні інформаційні ресурси є критичними каналами комунікації. Надмірна кількість спам-повідомлень у таких системах може призводити до перевантаження інфраструктури, зниження оперативності обробки звернень та підвищення ризиків несанкціонованого доступу до службової інформації.

У межах дослідження розглянуто підходи до попередньої обробки текстових повідомлень, включаючи токенізацію, нормалізацію, видалення стоп-слів та формування числових ознак за допомогою TF-IDF. Зазначені методи дозволяють ефективно перетворювати текстові дані у формат, придатний для використання алгоритмами машинного навчання. Для експериментального аналізу використано відкритий набір даних SMSSpamCollection, який містить реальні приклади спам- та не спам-повідомлень.

З метою досягнення поставленої мети здійснено порівняльний аналіз декількох методів класифікації, які найчастіше застосовуються у практичних задачах фільтрації спаму, а саме: наївного байесівського класифікатора, методу k-ближчих сусідів та дерева рішень. Порівняння методів проводилося за єдиним набором вхідних даних та однаковими умовами навчання, що забезпечує коректність отриманих результатів.

Практична реалізація моделей виконувалася в середовищі MATLAB із використанням вбудованих функцій для навчання та оцінювання класифікаторів. Для кожного методу обчислено основні метрики якості, зокрема Accuracy, Precision, Recall та F1-score, а також проаналізовано часові витрати на навчання і тестування моделей. Додатково проведено оцінювання якості класифікації за допомогою ROC-кривих та показника AUC.

Порівняльний підхід дозволив виявити сильні та слабкі сторони кожного з методів з точки зору їх практичного використання в інформаційних системах

державних органів, де важливими є не лише точність класифікації, а й швидкодія, масштабованість та стабільність роботи алгоритмів.

Результати порівняльного аналізу засвідчили, що наївний байєсівський класифікатор забезпечує найкраще співвідношення точності, швидкодії та стабільності результатів, що робить його доцільним для використання у реальних системах фільтрації спаму. Метод k-ближчих сусідів продемонстрував конкурентну точність, однак виявився менш ефективним за часовими характеристиками. Дерево рішень, у свою чергу, показало задовільні результати, але потребує додаткового налаштування для зменшення ризику перенавчання.

Отримані результати свідчать про доцільність використання порівняльного аналізу методів класифікації при виборі оптимальних алгоритмів для впровадження в інформаційні системи органів державної влади з урахуванням специфіки їх функціонування та вимог до інформаційної безпеки. Підтверджується ефективність застосування методів машинного навчання для задач автоматичної класифікації текстів і можуть бути використані при розробці інтелектуальних систем аналізу повідомлень у поштових сервісах, месенджерах та платформах цифрової комунікації.

Список використаних джерел

1. Jurafsky D., Martin J. H. Speech and Language Processing. Stanford University, 2023.
2. Kaddoura S., Alrabae S. A systematic literature review on spam content detection. Journal of Big Data, 2022.
3. Budiman D. Email spam detection: a comparison of SVM and Naive Bayes. Journal of Science Research & Engineering, 2024.