

УДК 323.266:004.9

DOI: 10.25140/2411-5363-2020-4(22)-91-95

Володимир Базилевич, Марія Прибителько

**СИСТЕМА ВИЯВЛЕННЯ ФЕЙКОВИХ НОВИН
ЗА ДОПОМОГОЮ DATA SCIENCE**

Актуальність теми дослідження. Великий об'єм інформації, що став традиційним для інформаційного суспільства, створює нові виклики для людства. На заміну проблемі складного доступу до інформації, що була актуальною раніше, приходять нова проблема: структуризації та фільтрування інформації. Серед загального потоку інформаційного шуму актуальним стає питання виокремлення правди. У цьому контексті важливим постає питання не лише розвитку критичного мислення, але й розробки технічних засобів виявлення фейків.

Постановка проблеми. Ця робота фокусується на понятті розробки системи виявлення фейкових новин, аналізі існуючих систем та їхніх принципів роботи, принципів побудови їхніх алгоритмів та особливостях їх використання.

Аналіз останніх досліджень та публікацій. Були розглянуті останні публікації у відкритому доступі, статистичні дані, звіти корпорації.

Виділення недосліджених частин загальної проблеми. Аналіз файлу буде виконаний за допомогою трьох методів/класифікаторів і без використання PassiveAgressive класифікатора. Обчислення та виведення результатів виконується за допомогою побудови матриць помилок та розрахування точності.

Постановка задачі. Основною метою роботи є створення на основі розглянутих матеріалів систему виявлення фейкових новин та досягти найбільш можливої точності.

Виклад основного матеріалу. Обрано вхідні дані для дослідження, проведена їх підготовка та аналіз. Проведено дослідження даних за допомогою методів/класифікаторів Логістичної регресії, Дерева рішень та Рандомного лісу. Обчислена точність виявлення фейкових новин.

Висновки відповідно до статті. Запропонована система дозволяє класифікувати новини як «фейкові» або «правдиві» з точністю 98-99 %.

Ключові слова: фейк; фейкова новина; правдива новина; класифікатор.

Рис.: 6. Табл. 1. Бібл.: 10.

Актуальність теми дослідження. На сьогодні актуальним є завдання аналізу правдивості інформації в новинах, якими заповнені всі існуючі канали отримання інформації. Її актуальність пов'язана з необхідністю запобігання паніки через отримання недостовірної інформації, розвінчування псевдонаукових фактів, що можуть загрожувати життю людей, боротьба з політичною пропагандою та інші.

Постановка проблеми. Відтоді як термін «фейкові новини» потрапив у повне використання та став загальноновживаним у 2016 році, ми стали свідками раціонального поширення неправдивої інформації та способів її поширення. У цьому році глядачам і читачам довелося стикнутися з фальшивими відеороликами, що були подані як дійсно правдиві, блогами, написаними анонімними «тролями», а також сумнівними теоріями, у які повірили навіть у Білому домі США. Це ще не кажучи про численні приклади помилкових або завідомо оманливих даних у ЗМІ або виступах всесвітньовідомих політиків.

Аналіз останніх досліджень і публікацій. Дослідження цієї проблеми проводилися такими вченими та організаціями, як Filip Mishevski [1], data-flair [2] та інші.

Виділення недосліджених частин загальної проблеми. У цій статті аналіз файлу буде виконаний за допомогою трьох методів/класифікаторів і без використання PassiveAgressive класифікатора. Обчислення та виведення результатів виконується за допомогою побудови матриць помилок та розрахування точності.

Постановка задачі. У статті проаналізована підготовка файлу до виявлення фейкових новин, робота з даними файлу. Реалізовано алгоритм виявлення фейкових новин.

Виклад основного матеріалу. «Фейкова новина» визначається як така, що «повністю складена і сфабрикована для обману читача, з метою збільшення трафіку і прибутку. Такі новини можна охарактеризувати як елемент інформаційної містифікації або навмисне поширення неправдивих фактів в онлайн- та традиційних медіа або ЗМІ з метою введення споживача в оману або отримання фінансової чи політичної вигоди» [3].

Ця робота фокусується на понятті розробки системи виявлення фейкових новин, аналізі існуючих систем та їхніх принципів роботи, принципів побудови їхніх алгоритмів та особливостях їх використання.

Для дослідження фейкових новин та отримання якомога точнішого результату потрібно мати досить великий обсяг вихідних даних. Причому чим більшу кількість новин вдасться проаналізувати, тим більшу точність результатів можна буде отримати. Серед чималої кількості варіантів був обраний сайт <https://www.kaggle.com/>, який надає можливість знайти багато сетів даних для майже будь-яких аналітичних потреб. У нашому випадку був обраний датасет «Fake and real news dataset» [4]. Цей набір даних складається із двох файлів Fake.csv та True.csv, які містять відповідно набір фейкових та набір правдивих новин. У датасеті використані новини, що були опубліковані з 31 березня 2015 року по 19 лютого 2018 року та з 13 січня 2016 року по 31 грудня 2018 року для файлів Fake.csv та True.csv відповідно.

Для зручності подальшого аналізу даних потрібно підготувати файл таким чином:

1. Додаємо лейбли “fake” та “true” для файлів.
2. Зливаємо обидва файли в один фінальний.
3. Перемішуємо дані, щоб отримати випадковий порядок новин у файлі.
4. Прибираємо колонки title та date, як такі, що не несуть цінності для дослідження.
5. Переводимо текст у lowercase.
6. Прибираємо пунктуацію.
7. Видаляємо “stopwords” (це такі слова, які не додають особливого значення реченню.

Вони можуть бути безпечно проігноровані без шкоди для змісту речення). Після виконання дій, описаних вище, файл можна вважати готовим для подальшого дослідження.

Перш ніж реалізувати алгоритм пошуку фейкових новин, проведемо первинний аналіз вмісту даних у файлі, адже важливо розуміти не лише з яким обсягом даних ми працюємо, а й з якими саме даними ми маємо справу.

Виведемо у вигляді стовпчастої діаграми кількість фейкових та правдивих новин. Для цього використаємо лейбли “true” та “false”, які були встановлені на початку роботи з файлом. Результат підрахунку представлений на рис. 1.

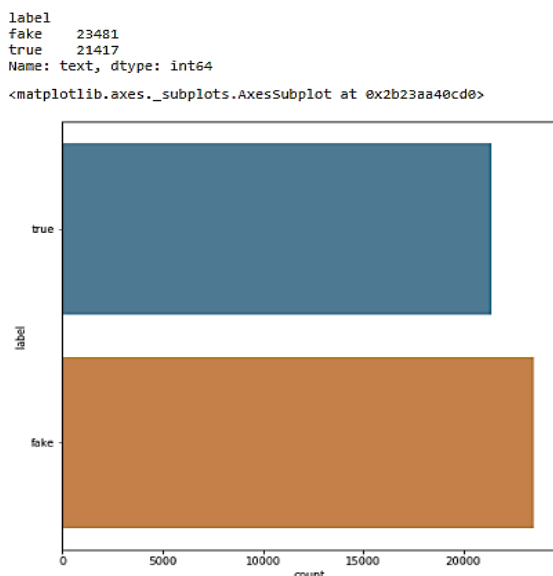


Рис. 1. Результат виведення кількості фейкових та правдивих новин

Виразуємо також точну кількість найбільш вживаних слів у фейкових та правдивих новинах і виведемо результат, як показано на рис. 2 та 3 відповідно.

Тепер можна зробити припущення про основну тему новин, що представлені у файлі, а також оцінити належність до певної географічної області (у нашому випадку країни) та опису дій певних груп населення (у нашому випадку можна побачити навіть конкретні прізвища).

Зробивши цей аналіз та розуміючи вміст файлу, можна переходити до створення алгоритму пошуку фейкових новин.

Процес моделювання складається з векторизації корпусу, що зберігається у стовпці «text», потім застосування TF-IDF (frequency-inverse document frequency) і, нарешті, алгоритму класифікації машинного навчання [5].

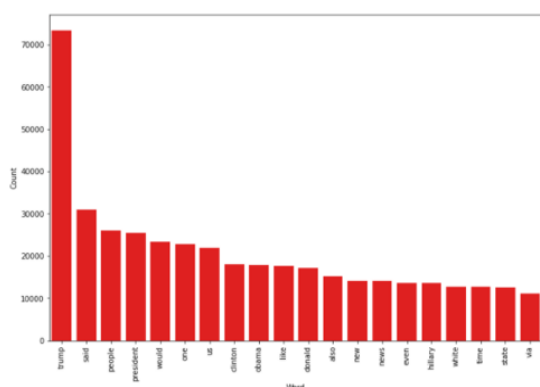


Рис. 2. Результат виведення кількості найбільш вживаних слів у фейкових новинах

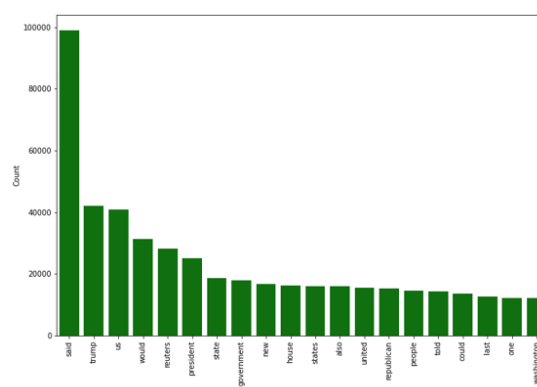


Рис. 3. Результат виведення кількості найбільш вживаних слів у правдивих новинах

Процес моделювання складається з векторизації корпусу, що зберігається у стовпці «text», потім застосування TF-IDF (frequency-inverse document frequency) і, нарешті, алгоритму класифікації машинного навчання [5].

У контексті цієї статті ми виділяємо три методи/класифікатори виявлення, а саме:

1. Логістична регресія.

Логістична регресія або (англ. *Logistic regression*) – це статистична модель, ої використовується для прогнозування ймовірності виникнення деякої події шляхом його порівняння з логістичною кривою. Ця регресія видає відповідь у вигляді ймовірності бінарної події (1 або 0) [6].

Основна ідея логістичної регресії полягає в тому, що простір вихідних значень може бути розділений лінійною границею (тобто прямою) на дві відповідних класам області [7].

2. Дерево рішень.

Дерево рішень – класифікатор, побудований на основі вирішальних правил виду «якщо, то», упорядкованих у деревоподібну ієрархічну структуру.

В основі роботи дерева рішень лежить процес рекурсивного розбиття вихідної безлічі об'єктів на підмножини, асоційовані з попередньо заданими класами. Розбиття проводиться за допомогою вирішальних правил, в яких здійснюється перевірка значень атрибутів за заданою умовою. Як навчальний набір даних використовується безліч спостережень, для яких попередньо задана мітка класу [8].

3. Рандомний(випадковий) ліс.

Випадковий ліс – один із найбільш приголомшливих алгоритмів машинного навчання, придуманий Лео Брейманом і Адель Катлер ще в минулому столітті. Він дійшов до нас у «первозданному вигляді» і є одним із небагатьох універсальних алгоритмів. Універсальність полягає, по-перше, в тому, що він хороший у багатьох завданнях, по-друге, в тому, що є випадкові ліса для вирішення завдань класифікації, регресії, кластеризації, пошуку аномалій, селекції ознак і т. ін. [9].

Побудуємо confusion matrix (матриці помилок) моделей для кожного із зазначених вище методів/класифікаторів [10]. На рис. 4 зображена матриця помилок для Логістичної регресії, на рис. 5 – для Дерева рішень та на рис. 6 – Рандомного лісу.

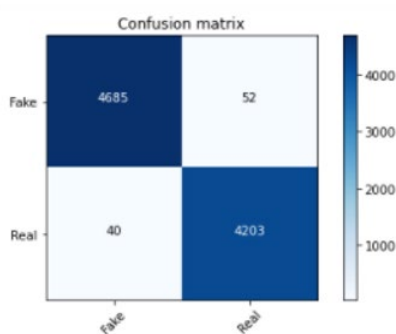


Рис. 4. Матриця помилок (Логістична регресія)

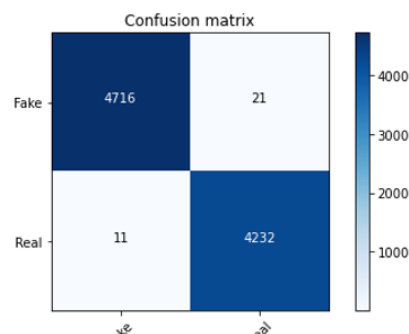


Рис. 5. Матриця помилок (Дерево рішень)

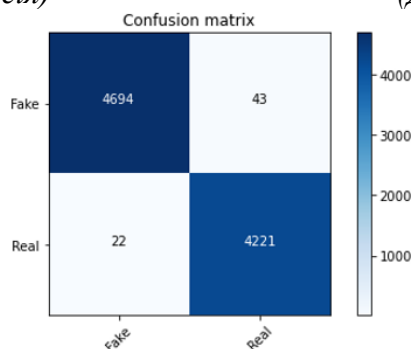


Рис. 6. Матриця помилок (Рандомний ліс)

Висновки відповідно до статті. Для зручності огляду, зробимо порівняльну таблицю отриманих результатів точності. Результат у вигляді можна побачити в табл. 1.

Таблиця 1

Порівняння результатів обчислення точності за різними методами

Назва методу	Отримана точність, %
Логістична регресія	98,98
Дерево рішень	99,64
Рандомний ліс	99,28

Отже, можемо зробити висновок, що результат виявлення фейкових новин за допомогою класифікатора Дерева рішень дав найбільший результат.

Список використаних джерел

- Mishevski F. Detecting Fake News With Python And Machine Learning. URL: <https://ichi.pro/ru/obnaruzhenie-fejkovyh-novostej-s-pomos-u-python-i-masinnogo-obucenia-272146875193862>.
- Detecting Fake News with Python. URL: <https://data-flair.training/blogs/advanced-python-project-detecting-fake-news>.
- Claire Wardle Fake news. It's complicated. URL: <https://firstdraftnews.org/latest/fake-news-complicated>.
- Fake and real news dataset. URL: <https://www.kaggle.com/clmentbisailon/fake-and-real-news-dataset>.
- TF-IDF. URL: <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>.
- Logistic regression. URL: https://en.wikipedia.org/wiki/Logistic_regression
- Saishruthi Swaminathan Logistic Regression – Detailed Overview. URL: <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>.
- Дерево рішень (Decision Trees) – Detailed Overview. URL: <https://wiki.loginom.ru/articles/decision-trees.html>.
- Random forest. URL: https://en.wikipedia.org/wiki/Random_forest.
- Confusion matrix. URL: https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html.

References

1. Mishevski F. Detecting Fake News With Python And Machine Learning. (n.d.). <https://ichi.pro/ru/obnaruzenie-fejkovyh-novostej-s-pomos-u-python-i-masinnogo-obucenia-272146875193862>.
2. Detecting Fake News with Python. (n.d.). <https://data-flair.training/blogs/advanced-python-project-detecting-fake-news>.
3. Claire Wardle Fake news. It's complicated. (n.d.). <https://firstdraftnews.org/latest/fake-news-complicated>.
4. Fake and real news dataset. (n.d.). <https://www.kaggle.com/clmentbisailon/fake-and-real-news-dataset>.
5. TF-IDF. (n.d.). <https://en.wikipedia.org/wiki/Tf%E2%80%93idf>.
6. Logistic regression. https://en.wikipedia.org/wiki/Logistic_regression
7. Saishruthi Swaminathan Logistic Regression – Detailed Overview. <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>.
8. Дерево рішень (Decision Trees) – Detailed Overview. <https://wiki.loginom.ru/articles/decision-trees.html>.
9. Random forest. https://en.wikipedia.org/wiki/Random_forest.
10. Confusion matrix. (n.d.). https://scikit-learn.org/stable/auto_examples/model_selection/plot_confusion_matrix.html.

UDC 323.266:004.9

Volodymyr Bazylevych, Maria Prybytko

FAKE NEWS DETECTION SYSTEM BASED ON DATA SCIENCE

Urgency of the research. Today, the task of analyzing the veracity of information in the news, which filled all existing channels for obtaining information, is relevant. Its urgency is related to the need to prevent panic by obtaining inaccurate information, debunking pseudo-scientific facts that can threaten people's lives, combating political propaganda and others.

Target setting This article focuses on the concept of developing a system for detecting fake news, analysis of existing systems and their principles of operation, principles of construction of their algorithms and features of their use.

Actual scientific researches and issues analysis. Recent open publications, statistics, and corporate reports were reviewed.

Uninvestigated parts of general matters defining. File analysis will be performed using three methods / classifiers and without the use of PassiveAggressive classifier. The calculation and derivation of results is performed by constructing error matrices and calculating accuracy.

The research objective. The main purpose of the work is to create a system for detecting fake news on the basis of the considered materials and to achieve the highest possible accuracy.

Presenting main material. Input data for the study were selected, prepared and analyzed. Data were studied using the methods / classifiers of Logistic Regression, Decision Tree and Random Forest. The accuracy of detecting fake news is calculated.

Conclusions. The proposed system allows to classify news as "fake" or "true" with an accuracy of 98-99 %.

Keywords: Python; fake; fake news; truthful news; classifier.

Fig.: 6. Table: 1. References: 10.

Базилевич Володимир Маркович – кандидат економічних наук, доцент, завідувач кафедри інформаційних та комп'ютерних систем, Національний університет «Чернігівська політехніка» (вул. Шевченка, 95, м. Чернігів, 14035, Україна).

Bazylevych Volodymyr – PhD in Economics, Associate Professor, Head of the Department of Information and Computer Systems, Chernihiv Polytechnic National University (95 Shevchenka Str., 14035 Chernihiv, Ukraine).

E-mail: bazvlamar@gmail.com

ORCID: <http://orcid.org/0000-0001-8935-446X>

ResearcherID: G-5764-2014

Scopus Author ID: 57193029322

Прибитько Марія Дмитрівна – здобувачка вищої освіти, Національний університет «Чернігівська політехніка» (вул. Шевченка, 95, м. Чернігів, 14035, Україна).

Prybytko Maria – PhD student, Chernihiv Polytechnic National University (95 Shevchenka Str., 14035 Chernihiv, Ukraine).

E-mail: prybytko.maria@gmail.com