

К.О. Бондалєтов¹
В.Б. Мокін¹
М.В. Григорчук¹
С.В. Джура¹
М.О. Кищук¹
О.В. Нерещуцький¹
С.Д. Неволя¹
А.М. Фурман¹
В.В. Гіжевський¹

ПОБУДОВА ДАТАСЕТУ ДЛЯ ТРЕНУВАННЯ ІНТЕЛЕКТУАЛЬНИХ МОДЕЛЕЙ ВЕБ-СИСТЕМИ З ІНФОРМАЦІЄЮ ПРО ЕКОЛОГІЧНІ ПРОБЛЕМИ ТА ЗАХОДИ У МАСИВАХ ВОД БАСЕЙНУ Р. ПІВДЕННИЙ БУГ WISEST-SBB

¹Вінницький національний технічний університет, Україна

Abstract

Робота присвячена побудові датасету для тренування інтелектуальних моделей веб-системи з інформацією про екологічні проблеми та заходи у масивах вод басейну р. Південний Буг «WISEST-SBB». Наведено розроблений шаблон, який наповнювали автори з використанням автоматичних та напіваавтоматичних способів веб-скрапінгу, парсингу та класифікації даних. Охарактеризовано результати випробування на цьому датасеті ряду інтелектуальних засобів для забезпечення автоматичності цих процесів у подальшому.

Ключові слова: NLP, Python, геоприв'язування даних, україномовний природномовний текст, штучний інтелект, масив вод, р. Південний Буг.

Abstract

The work is devoted to the construction of a dataset for training intelligent models of the web system with information about environmental problems and measures in the water bodies of the Southern Bug River basin "WISEST-SBB". The developed template, which was filled by the authors using automatic and semi-automatic methods of web scraping, parsing and data classification, is given. The results of the test on this dataset of a number of intelligent tools to ensure the automaticity of these processes in the future are characterized.

Keywords: NLP, Python, data geobinding, Ukrainian natural language text, artificial intelligence, water body, Southern Bug River.

Вступ

За рішенням басейнової ради Південного Бугу (протокол № 12 від 7 грудня 2022 року) на кафедрі системного аналізу та інформаційних технологій Вінницького національного технічного університету командою під керівництвом професора Мокіна В.Б. ведеться створення інтелектуальної веб-системи з інформацією про екологічні проблеми, природоохоронні заходи тощо про масиви вод басейну р. Південний Буг «WISEST Southern Bug Basin» (скорочено – «WISEST-SBB» або «WISESTR-SBB» – англ.: Water Information SystEm with Spatial and Temporal references for the Southern Bug Basin – «Водна інформаційна система з просторовою і часовою прив'язкою для басейну Південного Бугу») (<http://wisestr.ai/>). Така система дозволить здійснювати підтримку прийняття рішень під час розроблення та супроводу плану управління річковим басейном Південного Бугу [1, 2].

Наповнення системи здійснюється в напівавтоматичний спосіб, коли частина інформації знаходиться вручну, потім по ній тренуються інтелектуальні засоби для її пошуку та розмітки. Далі вручну перевіряється результат та удосконалюється. Потім – знов автоматичний етап і так – в циклі. Запорукою успіху є створення базового датасету для тренування цих інтелектуальних засобів.

Метою даного дослідження було створення датасету для тренування інтелектуальних моделей веб-системи з інформацією про екологічні проблеми та заходи у масивах вод басейну р. Південний Буг WISEST-SBB.

Систематизація інформації

Для розв’язання поставленої задачі Бондалетовим К.О. та Мокіним В.Б. було знайдено 20 перших фактів про стан масивів вод і на них відпрацьовано, по-перше, шаблон для їх систематизації, а по-друге, інтелектуальні засоби для їх оброблення з використанням зібраного раніше досвіду [3-5].

Було прийнято рішення роботи шаблон на базі MS Excel (рис. 1).

id	text	link	source	source type	geonames_text	geodata_text_other	date_doc	dates_text	dates	dates_text_other	status	issues	measures	comments
2	Екологи припустили погану якість Південного Бугу гинуть риби. Хмельницькі екологи перевірили рівень кисню та пошукували з інфекційними рибками. Кажуть, риба точно гине не через брак кисню. Доводяться яку причину морю назвали браком управління екологією та що кажуть хмельничани. Велику рибку пугача догори донизу опублікували в соцмережах. Напевно, мовля що в Південному Бугу відбуваються проблеми з екологією. В області центри контролю та профілактики хвороб повідомили про результати лабораторних досліджень. В управлінні з питань екології та контролю за біосферою, кажуть, також провели всі необхідні перевірки стану водних та прибережних смуг. Висновок, стверджують, що кисню в Південному Бугу вистачає. — Після останньої нічної перевірки кисню у річці, майже на рівні 0,7 мг/літр в усіх місцях проб (включено до уваги вант різноманітних розмірів) показав стабільні на рівні 4,0 мг/літр (додаток), хочемо констатувати, що риба	https://www.uv.gov.ua/rozdilnitskiy-otkaz-otkrytyi-sprava-s-otkazom-na-rybku-1163825.html	Веб-сторінка новин	сайт	Південний Буг, м. Хмельницький		28.07.2022		28.07.2022		1	1	1	
3	У Миколаївській області відбувається масовий мор риби: зникає частина русла Південного Бугу. Мисська влада та жителі вже почали вирішувати проблему. У місті Первомайська Миколаївської області на річці Південний Буг розгорілася екологічна катастрофа. Відбувається масовий мор риби. За її словами, через ремонт насосної станції за кілька днів був повністю зникла частина русла річки. «Що з рибкою добре видно на фото, люди зазначають без води (з мого в мого) і районі і без того самі проблеми», — написала Миколаївка. Рятує рибу місцевий житель Валерій Чечованов, ситуація спокійна дуже критична. Залишилася велика кількість риби, риба: різні породи риби в річці тоді, що є екологічною катастрофою, додалась ця проблема з тем, що риба річка вода впадає ще на 20 см, то ніяк нічого просто залишилось без води. Єдиний рішенням цього питання є відкрити річку води, але на це треба більше по темі не можуть надати нам порадювати кількість води і навіть, якщо вони відкрити не в воду, то в крайньому випадку вона буде йти	https://pnp.ukrainy.ua/mikolajivskiy-oblast-vidbuvaetsya-masovyj-mor-rybi-1335900.html	Веб-сторінка новин	сайт	Південний Буг, м. Первомайськ, Миколаївська область		30.09.2020		30.09.2020		0	1	0	

Рисунок 1. Структура шаблону з розміткою датасету

Як видно на рис. 1, шаблон передбачає інформацію про сам текст, його джерело, прив’язку у просторі (до карти) та у часі та класифікацію за рядом ознак. Ще можливими є примітки, коли розробник датасету не впевнений в деяких своїх оцінках цих даних.

Для розв’язання поставленої задачі Бондалетовим К.О. та Мокіним В.Б. було знайдено 20 перших фактів про стан масивів вод і на них відпрацьовано, по-перше, шаблон для їх систематизації, а по-друге, інтелектуальні засоби для їх оброблення з використанням зібраного раніше досвіду [3-5].

Тестування інтелектуальних засобів

Далі було перевірено яким чином розроблені Мокіним В.Б. та Бондалетовим К.О. інтелектуальні засоби зможуть здійснити геоприв’язку з використанням бібліотеки SpaCy та класифікацію знайденої інформації з використанням мовних моделей BERT та різних моделей машинного навчання для класифікації даних (рис. 2).

Score = 0.333 for WB code UA_M5.4_0002 and maybe: UA_M5.4_0002 UA_M5.4_0003 UA_M5.4_0004 UA_M5.4_0011 UA_M5.4_0013 UA_M5.4_0019 UA_M5.4_0028 UA_M5.4_0969 UA_M5.4_0970

- ['Південний Буг', 'Олександрівське водосховище'] 0.16666666666666663
- ['Південний Буг', 'Олександрівське водосховище'] 0.21428571428571422
- ['Південний Буг', 'Олександрівське водосховище'] 0.111111111111111109
- ['Південний Буг', 'Олександрівське водосховище'] 0.099999999999999996
- ['Південний Буг', 'Олександрівське водосховище'] 0.03448275862068965
- ['Південний Буг', 'Олександрівське водосховище'] 0.0
- ['Південний Буг', 'Олександрівське водосховище'] 0.0468749999999999986
- ['Південний Буг', 'Олександрівське водосховище'] 0.0
- ['Південний Буг', 'Олександрівське водосховище'] 0.06730769230769229
- ['Південний Буг', 'Олександрівське водосховище'] 0.08535533905932736
- ['Південний Буг', 'Олександрівське водосховище'] 0.07142857142857141

Рисунок 2. Приклад роботи NLP-програми на Python з автоматичним прив’язуванням інформації до масивів вод

Аналіз показав, що потрібно удосконалювати інтелектуальні засоби, враховуючи певні особливості української мови та певну стилістику подання новин, щоб підвищити точність автоматичної геоприв'язки та класифікації новин та іншої інформації.

Висновки

Робота присвячена побудові датасету для тренування інтелектуальних моделей веб-системи з інформацією про екологічні проблеми та заходи у масивах вод басейну р. Південний Буг «WISEST-SBB». Наведено розроблений шаблон, який наповнювали автори з використанням автоматичних та напівавтоматичних способів веб-скрапінгу, парсингу та класифікації даних. Охарактеризовано результати випробування на цього датасеті ряду інтелектуальних засобів для забезпечення автоматичності цих процесів у подальшому.

СПИСОК ВИКОРИСТАНОЇ ЛІТЕРАТУРИ

1. Водний кодекс України : Кодекс України від 06.06.1995 р. № 213/95-ВР : станом на 19 серп. 2022 р. URL: <https://zakon.rada.gov.ua/laws/show/213/95-вр#Text> (дата звернення: 07.06.2023).
2. Про затвердження Порядку розроблення плану управління річковим басейном : Постанова Каб. міністрів України від 18.05.2017 р. № 336. URL: <https://www.kmu.gov.ua/npas/249999756> (дата звернення: 04.06.2023).
3. Мокін В. Б. Інформаційна інтелектуальна технологія автоматизованої геоприв'язки екологічної текстової природно-мовної інформації / В. Б. Мокін, М. А. Гораш, Є. М. Крижановський, Т. Є. Вуж // Наукові праці ВНТУ [Електронний ресурс]. — 2020. — № 4. — Режим доступу: <https://praci.vntu.edu.ua/index.php/praci/article/view/624>
4. Vitalii Mokin. NLP for UA : BERT CLS & 10 Classifiers. Kaggle: Your Machine Learning and Data Science Community. URL: <https://www.kaggle.com/code/vbmokin/nlp-for-ua-bert-cls-10-classifiers> (date of access: 07.06.2023).
5. Vitalii Mokin. Kostiantyn Bondaletov. SpaCy for ukrainian text similarity. Kaggle: Your Machine Learning and Data Science Community. <https://www.kaggle.com/code/bondaletov/spacy-for-ukrainian-text-similarity> (date of access: 07.06.2023).

Бондалетов Костянтин Олегович — аспірант кафедри системного аналізу та інформаційних технологій; e-mail: bondaletov.k@gmail.com;

Мокін Віталій Борисович — д-р техн. наук, професор, завідувач кафедри системного аналізу та інформаційних технологій; e-mail: vbmokin@vntu.edu.ua;

Марина Василівна Григорчук — студент кафедри системного аналізу та інформаційних технологій; hruhorchuk08@gmail.com;

Сергій Вікторович Джура — студент кафедри системного аналізу та інформаційних технологій; dzurasergij4@gmail.com;

Максим Олександрович Кищук — студент кафедри системного аналізу та інформаційних технологій; makskishuk.78@gmail.com;

Олег Віталійович Неревський — студент кафедри системного аналізу та інформаційних технологій; onereutsk@gmail.com;

Сергій Дмитрович Неволя — студент кафедри системного аналізу та інформаційних технологій; nevolya2003@gmail.com;

Анна Михайлівна Фурман — студент кафедри системного аналізу та інформаційних технологій; furmanania1@gmail.com;

Владислав Віталійович Гіжевський - студент кафедри системного аналізу та інформаційних технологій; vladgiz2000@gmail.com.

Bondaletov Kostiantyn O. — Post-graduate student of the Chair of System Analysis and Information Technology, e-mail: bondaletov.k@gmail.com;

Mokin Vitalii B. — Dr. Sc. (Eng.), Professor, Head of the Chair of System Analysis and Information Technology, e-mail: vbmokin@vntu.edu.ua;

Hryhorchuk Maryna V. — student of Vinnytsia National Technical University, e-mail: hruhorchuk08@gmail.com;

Dzhura Serhii V. — student of Vinnytsia National Technical University, e-mail: dzurasergij4@gmail.com;

Kyshchuk Maksym O. — student of Vinnytsia National Technical University, e-mail: makskishuk.78@gmail.com;

Nereutskiy Oleh V. — student of Vinnytsia National Technical University, e-mail: onereutsk@gmail.com;

Nevolia Serhii D. — student of Vinnytsia National Technical University, e-mail: nevolya2003@gmail.com;

Furman Anna M. — student of Vinnytsia National Technical University, e-mail: furmanania1@gmail.com.

Hizhevskiy Vladyslav V. — student of Vinnytsia National Technical University, e-mail: vladgiz2000@gmail.com.