

**Лупей Максим Іванович**Аспірант кафедри інформаційних управляючих систем та технологій, [orcid.org/0000-0003-3440-3919](https://orcid.org/0000-0003-3440-3919)

Ужгородський національний університет, Ужгород

**ВИЗНАЧЕННЯ СТИЛЬОВОЇ НАЛЕЖНОСТІ ТЕКСТУ  
ЗА ДОПОМОГОЮ ШТУЧНИХ НЕЙРОННИХ МЕРЕЖ**

***Анотація.** Досліджено проблему розроблення ефективного способу визначення стильової належності текстів. Розглянуто такі стилі, як науковий, публіцистичний та офіційно-діловий. Для аналізу відібрані тексти одної тематики – про мову. Розглянуто різні поєднання методів векторизації та архітектур штучних нейронних мереж, які б забезпечили високий рівень розпізнаваності. Серед архітектур штучних нейронних мереж розглянуто: Support Vector Machines (SVM) (C-Support Vector Classification (SVC), Epsilon-Support Vector Regression (SVR)) та Multi Layer Perceptron (MLP). Серед методів векторизації розглянуто: HashingVectorizer, CountVectorizer та TfidfVectorizer. Проведені дослідження засвідчили, що всі розглядувані підходи найбільш ефективно розрізняють офіційно-ділові тексти, що пояснюється їх найбільшою стандартизованістю. Особливо ефективно розрізняються науковий та офіційно-діловий стилі. Найменшу точність розглядуваних методи показують при визначенні стильової приналежності, коли одним зі стилів є публіцистичний. Найбільш ефективним підходом для визначення стильової приналежності виявилось поєднання методу векторизації tfidfVectorizer та обох архітектур штучних нейронних мереж Support Vector Machines. На попередньому етапі для збільшення ефективності використовувався стемінг слів. У текстах, що містять не менше 500 символів, такий підхід допоміг забезпечити точність 94–98%, а час для навчання штучної нейронної мережі при цьому не перевищує одну секунду на комп'ютерах стандартної на цей час конфігурації. За допомогою бібліотеки Lime наведено візуалізацію дослідження роботи штучної нейронної мережі, що є надзвичайно важливим емпіричним матеріалом для фахівців-філологів для проведення подальшого лінгвістичного аналізу.*

**Ключові слова:** стиль; класифікація; корпусна лінгвістика; штучні нейронні мережі; векторизація тексту

**Вступ**

Стрімкий розвиток технологій машинного навчання уможливує розгляд традиційних наукових питань під новим кутом зору. Одним із цікавих і перспективних напрямів є аналіз особливостей природної мови, зокрема й української, за допомогою штучних нейронних мереж (ШНМ). Дослідження значних масивів даних, робота з великими корпусами текстів з використанням комп'ютерних технологій може не тільки спростити лінгвістам рутинні процедури (розписування текстів на картки, вибіркового аналізу мовних одиниць), а й допомогти розвинути нові аспекти й напрями лінгвістичного аналізу.

Кількість текстів, написаних українською мовою, невпинно зростає. Значна частина текстів представлена онлайн, що спрощує доступ до емпіричної бази досліджень. Ці тексти мають різне функціональне призначення, у них використовуються мовні засоби, що відповідають меті й завданням спілкування у певній комунікативній ситуації.

**Аналіз літературних даних  
та постановка проблеми**

Традиційно в українському мовознавстві виокремлюють розмовний, публіцистичний, науковий, офіційно-діловий, художній і конфесійний функціональні стилі, однак деякі дослідники виокремлюють також епістолярний та ораторський. Диференціація функціональних стилів здійснюється на основі власне мовних ознак (кількісні та якісні показники мовно-структурних одиниць) та екстралінгвальних чинників (суспільна сфера функціонування текстів). Відома українська лінгвістка С. Я. Єрмоленко зауважує, що «релевантними для функціональних стилів є сфера суспільної діяльності, тип мислення, стереотипні конструкти, тобто призначення, мета текстів, кількість використовуваних одиниць і їх стильове навантаження». З огляду на наведені чинники лінгвісти виокремлюють типові ознаки кожного функціонального стилю української мови [1].

У роботі [2] розроблено метод на мові програмування Python, за допомогою якого

проведено аналіз приналежності тексту до наукового, художнього, офіційно-ділового та публіцистичного. Проаналізовано за допомогою цього інструменту по 65 текстів із зібраного корпусу, було класифіковано 88% із них. Причому найбільша точність була досягнута при класифікації ділових та художніх текстів. Відмітимо, що найбільша похибка була при розмежуванні публіцистичних та наукових текстів.

У роботі [3] розглядається бібліотека для машинного навчання Scikit-learn мови програмування Python. За допомогою цієї бібліотеки проведено аналіз приналежності текстів до категорії: релігія (християнство), медицина, атеїзм та графіка. Досягнута класифікатором точність становить 91%.

Загалом можна стверджувати, що проблематика визначення стильової належності тексту потребує подальших досліджень. Наявність величезної кількості різних інструментів дає можливість отримати більш високі результати, звуживши досліджувану область. Тому у цій роботі було розглянуто лише три стилі: науковий, публіцистичний та офіційно-діловий, а також новий стек сучасних технологій.

### Мета статті

Метою цієї статті є визначення належності текстів до наукового, публіцистичного та офіційно-ділового стилів за допомогою ШНМ.

Для реалізації мети поставлено такі завдання:

- добрати тексти, що належать до наукового, публіцистичного та офіційно-ділового стилів і репрезентують близьку тематику;
- розробити програму для парсингу текстів і підготовки для роботи з ними ШНМ;
- визначити найбільш ефективні методи векторизації текстів;
- визначити ефективну структуру ШНМ для визначення належності текстів до досліджуваних функціональних стилів української мови.

### Виклад основного матеріалу

#### Метод дослідження

Оскільки щодо текстів застосовується процедура векторизації, яка полягає у перетворенні колекцій текстових документів у частотні матриці лексем, то основними одиницями, на основі яких ШНМ визначає належність тексту до певного стилю, є лексеми й афікси.

З огляду на це окреслимо основні особливості досліджуваних функціональних стилів на лексичному й словотворочому рівнях [1]:

- найвища міра книжності, використання стандартизованих висловів, штампів, кліше, книжних зворотів тощо, висока частотність

вживання абстрактних іменників на *-ість*, *-ання* з термінологічним значенням – в офіційно-діловому стилі;

- використання термінології, загальнонавчаних слів тільки в одному з притаманних їм значень (функціональна однозначність слів), використання абстрактних іменників із суфіксами *-ість*, *-ств(о)*, *-от(а)*, *-аці(я)*, *-изм* тощо – у науковому стилі;

– використання засобів увиразнення мовлення – стилістичних фігур, образності, інверсії та ін. – з метою емоційного впливу на реципієнта, вживання лексичних синонімів для надання висловом «небуденності», використання словотворчих форм із префіксом *не-*, паралельне вживання книжних і розмовних лінгвальних одиниць – у публіцистичному стилі.

Отже, в офіційно-діловому, науковому й публіцистичному стилях різною мірою може виражатися авторське «я» (індивідуальні особливості мовлення автора, його мовні вподобання): найменше воно виражається в офіційно-діловому, найбільше – у публіцистичному стилі. Для офіційно-ділового й наукового стилів не характерне вживання таких стилістично маркованих лексичних одиниць, як діалектизми, розмовні слова, емоційна й експресивна лексика тощо.

Вибір для дослідження текстів саме цих трьох функціональних стилів зумовлений, з одного боку, досить високим рівнем «книжності» й стандартизації офіційно-ділового й наукового стилів, з іншого – чіткістю, логічністю й одночасним емоційно-експресивним спрямуванням публіцистики. Оскільки, відповідно до методики нашого дослідження, штучній нейронній мережі надається тільки корпус текстів без вказівки їх стильової належності й без характеристики типових ознак певного стилю, то цікавим є те, який рівень точності результатів буде для кожного зі стилів.

### Результати дослідження

Для аналізу відібрані тексти однієї тематики – про мову. Тексти однієї тематики вибрано з метою мінімізації ймовірності правильного визначення належності тексту до певної групи за тематичною ознакою й актуалізації власне лінгвальних критеріїв належності тексту до того чи іншого функціонального стилю.

Офіційно-ділові тексти відібрані із сайту Верховної Ради України ([zakon.rada.gov.ua](http://zakon.rada.gov.ua)): Європейська хартія регіональних мов або мов меншин, закони «Про мови в Українській РСР», «Про забезпечення функціонування української мови як державної», а також Укази Президента України, рішення Конституційного Суду України, листи Міністерства освіти і науки України з мовної проблематики. Наукові тексти взяті з видання

«Мовні права в сучасному світі (Збірник наукових праць)» (Ужгород, 2014. – 351 с.). Щодо наукових текстів зроблено попереднє опрацювання для того, щоб залишити власне авторські тексти (вилучено бібліографічні посилання до статей, а також додатки, у яких уміщені документи).

Публіцистичні тексти з мовної проблематики дібрано з видань «Український тиждень», «Дзеркало тижня», «Українська правда», «Українська літературна газета». Авторами статей є такі публіцисти, як Оксана Форостина, Олександр Рудяченко, Дарія Гайдай, Аскольд С. Лозинський, Віктор Грабовський, Володимир Михайленко, Ярина Бусол, Павло Зуб'юк, Світлана Єременко та ін.

Кількісні показники відібраних текстів такі:

- наукові тексти – близько 81 тисяча слововживань (619 тисяч символів);
- офіційно-ділові тексти – 32 тисячі слововживань, 248 тисяч символів;
- публіцистичні тексти – 16 тисяч слововживань (116 тисяч символів).

Загальна схема запропонованого методу наведена на рис. 1.



Рисунок 1 – Загальна схема роботи методу

Задачі, які належать до text mining ефективно можуть бути розв'язані за допомогою ШНМ [4]. Саме ШНМ лягли в основу цього дослідження.

Як і в роботі [5], досліджувалась ефективність різних поєднань методів векторизації та ШНМ (рис. 1). Були розглянуті такі методи векторизації: HeshingVectorizer, CountVectorizer та TfidVectorizer. Серед ШНМ було розглянуто: Support Vector Machines (SVM) (C-Support Vector Classification (SVC), Epsilon-Support Vector Regression (SVR)) та Multi Layer Perseptron (MLP). У цій роботі на попередньому етапі було застосовано стемінг слів (процес скорочення слова до основи шляхом відкидання допоміжних частин, а в нашому випадку

конкретно закінчень), що дало можливість покращити результат в середньому на близько 1%. Приклад тесту із застосуванням процедури стемінгу наведено в табл. 1.

Після проведення експерименту виявилось, що найбільш ефективні результати дає поєднання методу векторизації tfidfVectorizer та архітектур штучних нейронних мереж SVC та SVR (табл. 2). Попереднє опрацювання вхідного тексту методом векторизації TfidfVectorizer наведено в табл. 1. Використання методу векторизації HeshingVectorizer в поєднанні з розглядуваними ШНМ дає точність на 2 – 3% меншою, ніж метод векторизації TfidfVectorizer. Метод векторизації CountVectorizer в поєднанні з розглядуваними ШНМ показує найнижчі результати.

Як видно з табл. 2, найбільш ефективно з використанням описаних засобів можна розрізнити науковий та офіційно-діловий стилі. Всі три розглядувані ШНМ в поєднанні з методом векторизації tfidfVectorizer забезпечують точність понад 97%. Розрізнити науковий та публіцистичний стилі цими ж засобами вдається з точністю понад 93%. А ось виявити різницю між публіцистичним та офіційно-діловим за допомогою ШНМ MLP можна лише з точністю 84,54%, що суттєво менше, ніж за допомогою інших ШНМ.

Якщо порівнювати між собою ШНМ SVC та SVR, то незначну перевагу має перша. Підсилити цю перевагу допомагає наявна для неї бібліотека Lime, що дає змогу візуалізувати роботу з класифікації тексту (рис. 2 – 4). Це стає надзвичайно цінним емпіричним матеріалом подальших досліджень для філологів. Аналіз приналежності послівно (рис. 2), підсвітка розпізнавання (рис. 3) та ймовірносний розділ тексту між стилями (рис. 4) є важливими наочними засобами, які показують роботу ШНМ всередині, а забезпечена висока точність допомагає фахівцям-лінгвістам робити певні важливі узагальнення.

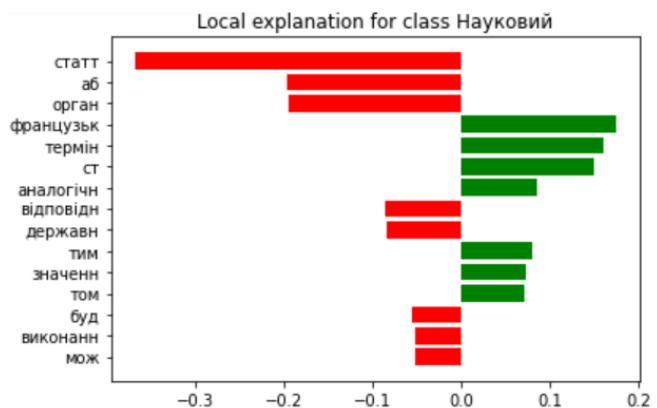


Рисунок 2 – Візуалізація роботи за допомогою бібліотеки Lime

статт 14 забороня юридичн особ уживанн виробнич, торгівельн аб сервіс марок на баз іноземн вислов аб термін, якщ аналогічн вислов аб термін з тим же значенн затверджен відповідн до ум, передбачен нормативн документ про збагаченн французьк мов. статт 15 вимаг від центральн і територіальн державн орган наданн буд-як субсид лиш за умов виконанн отримувач положен дан закон. невиконанн закон отримувач мож потягт за соб повн аб частк поверненн субсиді (субвенці). статт 16 опис компетенц представник правоохоронн орган, у том числ карн розшук, щод встановленн факт порушен норм, як впливають зі ст. ;

Рисунок 3 – Приклад розпізнавання офіційно-ділового стилю з датасету науковий / офіційно-діловий (науковий: жовто-оранжевий, офіційно-діловий: блакитно-синій)

Таблиця 1 – Приклад процедури стемінгу для конкретного тексту

Індекс в датасеті	Приклад тексту
Source	Введенням у конституційно-правове поле України як країни з поліетнічним складом населення припису про функціонування української мови у всіх сферах суспільного життя на всій території України конституцієдавець визначив її специфічну роль, зокрема, пов'язавши різноаспектну діяльність держави як суспільного регулятора з обов'язком послуговуватися уніфікованим мовним засобом суспільної комунікації – українською мовою. Одночасно це є демонстрацією ототожнення державного утворення України з титульною українською нацією. – Основний Закон України є програмним документом, якому надано юридичну форму, де цілі мають першорядне значення
Stemmed	введенн у конституційн-прав пол україн як країн з поліетнічн склад населенн припис про функціонуванн українськ мов у всіх сфер суспільн житт на всі територі україн конституцієдавець визнач і специфічн рол, зокрем, пов'яза різноаспектн діяльніст держав як суспільн регулятор з об'язк послуговув уніфікован мовн засоб суспільн комунікаці – українськ мов. одночасн це є демонстраціє ототожненн державн утворенн україн з титульн українськ націє. – основн закон україн є програмн документ, як надан юридичн форм, де ціл мают першорядн значенн.

Таблиця 2 – Результати експериментів при поєднанні методу векторизації tfidfVectorizer та архітектур штучних нейронних мереж

№	K-FOLD	Структура ШНМ	VALUES	ACCURACY	Тип
1	10/0.3	SVM (SVC)	F1 ≈ 0.9790	0.9680	Науковий Офіційно-діловий
2	10/0.3	SVM (SVC)	F1 ≈ 0.9454	0.9042	Науковий Публіцистичний
3	10/0.3	SVM (SVC)	F1 ≈ 0.9552	0.9654	Публіцистичний Офіційно-діловий
4	10/0.3	SVM (SVR)	F1 ≈ 0.9812	0.9711	Науковий Офіційно-діловий
5	10/0.3	SVM (SVR)	F1 ≈ 0.9406	0.8946	Науковий Публіцистичний
6	10/0.3	SVM (SVR)	F1 ≈ 0.9534	0.9639	Публіцистичний Офіційно-діловий
7	10/0.3	MLP	F1 ≈ 0.9774	0.9651	Науковий Офіційно-діловий
8	10/0.3	MLP	F1 ≈ 0.9305	0.8750	Науковий Публіцистичний
9	10/0.3	MLP	F1 ≈ 0.8454	0.8955	Публіцистичний Офіційно-діловий

Важливою особливістю реалізації цього методу є те, що не потрібно було порівнювати відразу три стилі, а достатньо перевірити два рази по два стилі і це покращувало точність та допомагало проводити

швидке навчання нейронної мережі. Таке рішення є досить простим, але дуже ефективним. Додатковими можливостями для покращення точності є використання засобів, описаних в [6; 7].

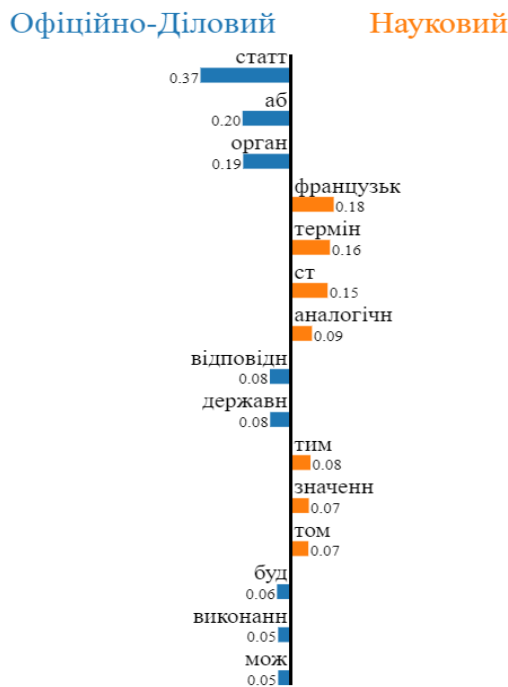


Рисунок 4 – Приклад роботи методу при визначенні належності тексту

У роботі проведено дослідження залежності точності визначення приналежності до стилю при поєднанні методу векторизації tfidfVectorizer та архітектури SVC від часу навчання ШНМ (рис. 5). Визначено, що достатньо однієї секунди для забезпечення високої точності на комп'ютері з процесором Intel Core i7 1.80 ГГц 32 ГБ ОЗУ.

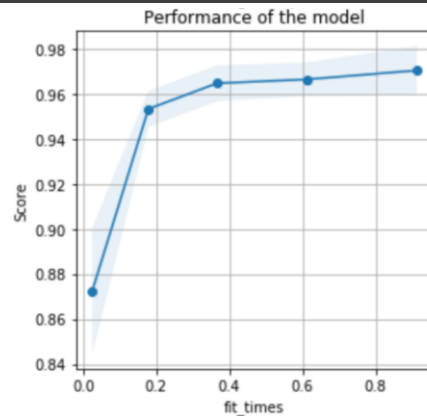


Рисунок 5 – Зміна точності при збільшенні часу в секундах

## Висновки

Проведене дослідження дало змогу виявити найбільш ефективне поєднання методу векторизації та ШНМ. Наявність візуалізації роботи розробленого інструменту є цікавим емпіричним матеріалом для подальших лінгвістичних студій, оскільки з погляду філології важливо з'ясувати, завдяки яким саме мовним елементам визначається стильова приналежність. Продовжуючи дослідження в обраному напрямі у подальшому, вважаємо за доцільне розширити джерельну базу дослідження: по-перше, використати більші корпуси текстів обраних стилів (наукового, публіцистичного, офіційно-ділового); по-друге, проаналізувати за допомогою запропонованої методики текстів, що належать до інших стилів (художній, конфесійний, розмовний).

## Список літератури

1. Єрмоленко, С.Я. Лінгвостилістика: основні поняття, напрями й методи дослідження. Українська лінгвостилістика XX – початку XXI ст.: система понять і бібліографічні джерела [уклад.: Бибик С.П., Єрмоленко С.Я., Коць Т.А. та ін.; за ред. д-ра філол. наук, проф. С.Я. Єрмоленко]. – К.: Грамота, 2007.
2. Дубовик, А.Р. Автоматическое определение стилистической принадлежности текстов по их статистическим параметрам. Компьютерная лингвистика и вычислительные онтологии 1 (2017): 29 – 45.
3. Pedregosa F. et al. "Scikit-learn: Machine learning in Python." the Journal of machine Learning research 12 (2011): 2825 – 2830.
4. Bodyanskiy, Y. "Computational Intelligence Techniques for Data Analysis" in Leipziger Informatik-Tage, 2005, pp. 15 – 36.
5. Lupei, M., Mitsa, A., Repariuk, V., & Sharkan, V., (2020). Identification of authorship of Ukrainian-language texts of journalistic style using neural networks. Eastern-European Journal of Enterprise Technologies, 1 (2 (103)), 30 – 36. doi: <https://doi.org/10.15587/1729-4061.2020.195041>
6. Bodyanskiy et al. "Deep 2D-Neural Network and its Fast Learning," in Proceedings of the 2018 IEEE Second International Conference on Data Stream Mining & Processing, DSMP 2018. Lviv, Ukraine, 21 – 25 August 2018, pp. 519 – 523.
7. Rashkevych, Y., Peleshko, D., and Pasyeka, M. "Optimization search process in database of learning system," in Proceedings of the 2nd IEEE International Workshop on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications, IDAACS 2003. Lviv, Ukraine, 8 – 10 Sept. 2003, pp. 358 – 361.

Стаття надійшла до редколегії 02.04.2020

Лупей Максим Іванович

Аспірант кафедри інформаційних управляючих систем і технологій, [orcid.org/0000-0003-3440-3919](https://orcid.org/0000-0003-3440-3919)  
Ужгородський національний університет, Ужгород

## ОПРЕДЕЛЕНИЕ СТИЛЕВОЙ ПРИНАДЛЕЖНОСТИ ТЕКСТА С ПОМОЩЬЮ ИСКУССТВЕННЫХ НЕЙРОННЫХ СЕТЕЙ

**Аннотація.** Исследована проблема разработки эффективного способа определения стилиевой принадлежности текстов. Рассмотрены такие стили: научный, публицистический и официально-деловой. Для анализа отобраны тексты одной тематики – о языке. Рассмотрены различные сочетания методов векторизации и архитектур искусственных

нейронных сетей, которые могут обеспечить высокий уровень узнаваемости. Среди архитектур искусственных нейронных сетей рассмотрены: Support Vector Machines (SVM) (C-Support Vector Classification (SVC), Epsilon-Support Vector Regression (SVR)) и Multi Layer Perceptron (MLP). Среди методов векторизации рассмотрены: HeshingVectorizer, CountVectorizer и TfidVectorizer. Проведенные исследования показали, что все рассматриваемые подходы наиболее эффективно различают официально-деловые тексты, что объясняется их большей стандартизованностью. Особенно эффективно различаются научный и официально-деловой стили. Наименьшую точность рассматриваемые методы показывают при определении стилиевой принадлежности, когда одним из стилей является публицистический. Наиболее эффективным подходом для определения стилиевой принадлежности оказалось сочетание метода векторизации tfidfVectorizer и двух архитектур искусственных нейронных сетей Support Vector Machines. На предварительном этапе для повышения эффективности использовался стемминг слов. В текстах, содержащих не менее 500 символов, такой подход позволил обеспечить точность 94 – 98%, а время для обучения искусственной нейронной сети при этом не превышает одну секунду на компьютерах стандартной конфигурации. С помощью библиотеки Lime приведена визуализация исследования работы искусственной нейронной сети, что является чрезвычайно важным эмпирическим материалом для специалистов-филологов при проведении дальнейшего лингвистического анализа.

**Ключевые слова:** стиль; классификация; корпусная лингвистика; искусственные нейронные сети; векторизация текст

**Lupei Maksym**

Postgraduate student, Department of Information Management Systems and Technologies, [orcid.org/0000-0003-3440-3919](https://orcid.org/0000-0003-3440-3919)  
Uzhgorod National University, Uzhgorod

#### DETERMINING THE STYLISTIC AFFILIATION OF THE TEXT USING ARTIFICIAL NEURAL NETWORKS

**Abstract.** The research is about the problem of developing an effective way to determine the stylistic affiliation of texts. Styles such as scientific, journalistic and official-business are considered. Texts of one subject – about language – selected for the analysis. Different combinations of vectorization methods and architectures of artificial neural networks considered which would provide a high level of recognition. Among the architectures of artificial neural networks – Support Vector Machines (SVM) (C-Support Vector Classification (SVC), Epsilon-Support Vector Regression (SVR)) and Multi Layer Perceptron (MLP). Among the vectorization methods – HeshingVectorizer, CountVectorizer and TfidVectorizer. Studies have shown that all the approaches considered most effectively distinguish between official and business texts, due to their greatest standardization. Scientific and official business styles are especially effective. The considered methods show the least accuracy in determining the stylistic affiliation, when one of the styles is journalistic. The most effective approach to determining stylistic affiliation was the combination of the vectorization method and both architectures of artificial neural networks Support Vector Machines. In the previous stage, word stemming was increasing efficiency. In texts with at least 500 characters, this approach has ensured an accuracy of 94-98%, and the time for learning an artificial neural network does not exceed one second on computers of the standard configuration at this time. With the help of the Lime library, a visualization of the study of the operation of an artificial neural network presents, which is extremely important empirical material for philologists for further linguistic analysis.

**Keywords:** style; classification; corpus linguistics; artificial neural networks; vectorization of the text

#### References

1. Ermolenko, S.Ya., (2007). *Linguostilistic: main terms, directions and methods of investigation. Ukrainian linguostilistic XX – beginning of XXI century: systems of terms and bibliography.* K.: Gramota.
2. Dubivik, A.R., (2017). Automatic determination of stylistic belonging of text at their statistic parameters. *Computer linguistics and calculation ontology*, 1, 29 – 45.
3. Pedregosa, F. et al. (2011). Scikit-learn: Machine learning in Python. *Journal of machine Learning research*, 12, 2825 – 2830.
4. Bodyanskiy, Y., (2005). *Computational Intelligence Techniques for Data Analysis.* Leipziger Informatik-Tage, 15 – 36.
5. Lupei, M., Mitsa, A., Repariuk, V., & Sharkan, V., (2020). Identification of authorship of Ukrainian-language texts of journalistic style using neural networks. *Eastern-European Journal of Enterprise Technologies*, 1 (2 (103)), 30 – 36. doi: <https://doi.org/10.15587/1729-4061.2020.195041>
6. Bodyanskiy, Y. et al. (2018). Deep 2D-Neural Network and its Fast Learning. *Proceedings of the 2018 IEEE Second International Conference on Data Stream Mining & Processing, DSMP 2018.* Lviv, Ukraine, 21 – 25 August 2018, pp. 519 – 523.
7. Rashkevych, Y., Peleshko, D., and Pasyeka, M., (2003). Optimization search process in database of learning system. *Proceedings of the 2nd IEEE International Workshop on Intelligent Data Acquisition and Advanced Computing Systems: Technology and Applications, IDAACS 2003.* Lviv, Ukraine, 8-10 Sept. 2003, pp. 358 – 361.

#### Посилання на публікацію

- APA Lupei, Maksym, (2020). Determining the stylistic affiliation of the text using artificial neural networks. *Management of Development of Complex Systems*, 42, 63 – 68, [dx.doi.org/10.32347/2412-9933.2020.42.63-68](https://doi.org/10.32347/2412-9933.2020.42.63-68).
- ДСТУ Лупей М.І. Визначення стильової належності тексту за допомогою штучних нейронних мереж [Текст] / М.І. Лупей // Управління розвитком складних систем. – 2020. – № 42. – С. 63 – 68, [dx.doi.org/10.32347/2412-9933.2020.42.63-68](https://doi.org/10.32347/2412-9933.2020.42.63-68).