

Отримано: 29 березня 2023 р.

Прорецензовано: 20 квітня 2023 р.

Прийнято до друку: 1 травня 2023 р.

e-mail: kalina.ua1980@gmail.com

ORCID ID: <https://orcid.org/0000-0002-5275-4715>

DOI: 10.25264/2519-2558-2023-17(85)-59-61

Калініченко М. М. Застосування семантичних моделей у літературознавстві та експертизі об'єктів авторського права. *Наукові записки Національного університету «Острозька академія»: серія «Філологія»*. Острог: Вид-во НаУОА, 2023. Вип. 17(85). С. 59–61.

УДК: 8.80

Калініченко Михайло Михайлович,
кандидат філологічних наук, доцент кафедри романо-германської філології,
Рівненський державний гуманітарний університет

ЗАСТОСУВАННЯ СЕМАНТИЧНИХ МОДЕЛЕЙ В ЛІТЕРАТУРОЗНАВСТВІ ТА ЕКСПЕРТИЗІ ОБ'ЄКТІВ АВТОРСЬКОГО ПРАВА

Стаття розглядає застосування семантичних моделей в галузях експертизи авторського права та літературного аналізу. У змісті представлено провідні типи семантичних моделей, що використовуються для аналізу текстів (LSA, NER, тематичне моделювання тем, розпізнавання іменованих сутностей та ін.). Хоча кожна модель має свої переваги та недоліки, існують істотні проблеми їхнього використання. Втім, семантичні моделі значно впливають на галузь, дозволяючи ефективніше та точніше виявляти випадки плагіату та порушення авторських прав. Огляд, представлений у публікації, стосовно сутності та особливостей використання семантичних моделей в експертизі авторського права та літературному аналізі дозволяє зробити наступні висновки. Семантичні моделі значно поліпшують ефективність аналізу тексту. Ці моделі використовують передові технології для виявлення структурних патернів, тем та схожих елементів у текстах, а також дозволяють експертам глибше зрозуміти твори, які вони аналізують.

У галузі експертизи авторського права семантичні моделі використовуються для аналізу значних обсягів тексту та виявлення можливих порушень авторських прав. Є кілька способів вирішення можливих проблем з використанням семантичних моделей, таких як комбінування кількісного та якісного аналізу, використання декількох моделей одночасно, врахування вибіркового спотворення аналітичної інформації, вивчення ширшого контексту та додаткова перевірка результатів експертами. Кожна семантична модель має власні переваги та обмеження. Отже, не існує жодної ідеальної моделі для усіх типів аналізу. Наприклад, застосування LSA призводить до помітних труднощів з ідентифікацією плагіату у тих випадках, коли порівнювані документи дуже відрізняються за змістом і структурою, тоді як тематичне моделювання потерпає від проблеми з ідентифікацією тонких семантичних структур в художніх текстах.

У кінцевому підсумку, використання семантичних моделей у галузі експертизи авторського права та літературного аналізу має велику практичну цінність, але важливо бути усвідомленим їхніх обмежень та використовувати їх у поєднанні з іншими методами аналізу для забезпечення найбільш об'єктивних та комплексних досліджень. З розвитком технологій будуть розроблені нові та більш вдосконалені семантичні моделі, що дозволять проводити ще більш детальний аналіз текстів.

Ключові слова: семантичні моделі, експертиза авторського права, літературний аналіз, Latent Semantic Analysis (LSA), тематичне моделювання, Named Entity Recognition (NER).

Mykhaylo Kalinichenko,
PhD, Associate Professor of the Department of Romance-Germanic Philology,
Rivne State University of the Humanities

APPLICATION OF SEMANTIC MODELS IN LITERARY STUDIES AND EXPERTISE OF COPYRIGHT OBJECTS

The article examines the application of semantic models in the fields of copyright expertise and literary analysis. The content showcases various types of semantic models utilized for text analysis, such as Latent Semantic Analysis (LSA), topic modeling, named entity recognition (NER), sentiment analysis, and dependency parsing. While each model has its own advantages and drawbacks, challenges in their utilization exist, such as over-reliance on quantitative analysis, model biases, lack of contextual information, and potential misinterpretation of meaning. Nevertheless, semantic models have significantly impacted the field by enabling the identification of instances of plagiarism and copyright infringement more efficiently and accurately. The overview presented in the publication on the essence and peculiarities of using semantic models in copyright expertise and literary analysis allows for the following conclusions. Semantic models significantly improve the efficiency of text analysis. These models use advanced technologies to identify structural patterns, themes, and similar elements in texts, as well as allowing experts to gain a deeper understanding of the works they are analyzing. In the field of copyright expertise, semantic models are used for analyzing significant volumes of text and identifying possible copyright infringements. Although semantic models have many advantages, they also have functional limitations. It is important to be aware of these limitations and use analytical models in combination with other methods to ensure a more comprehensive understanding of the text. There are several ways to address possible problems with using semantic models, including combining quantitative and qualitative analysis, using multiple models at once, taking into account selection biases of analytical information, studying the broader context, and additional verification of results by human experts. Ultimately, the use of semantic models in copyright expertise and literary analysis has enormous practical value, but it is important to be aware of their limitations and use them in conjunction with other analysis methods to ensure the most objective and comprehensive research. As technology develops, new and more advanced semantic models will be developed that will allow for even more detailed analysis of texts.

Keywords: semantic models, copyright expertise, literary analysis, Latent Semantic Analysis (LSA), topic modeling, named entity recognition (NER).

Постановка проблеми у загальному вигляді та зв'язок із важливими науковими чи практичними завданнями. У світі гуманітарної науки семантичні моделі стали популярними в 1970-х роках й незабаром революціонізували галузі літературознавства та комп'ютерної обробки текстів, дозволяючи машинам розуміти людську мову та виконувати завдання, такі як переклад і статистичні дослідження текстів. Використання семантичних моделей в галузі експертизи авторського права і літературного аналізу дозволяє аналізувати великі обсяги тексту та ідентифікувати можливі порушення авторських

прав. Семантичні моделі можуть виявляти схожості між текстами, навіть якщо вони мають різні фразеологію та граматику, аналізуючи імпліцитний семантичний зміст слів та їх комбінацій.

Відомо чимало резонансних справ, під час яких експерти та правники ефективно скористалися семантичними моделями для експертизи авторських прав. У 2004 році Google запустив проєкт з диджиталізації мільйонів книг з бібліотек по всьому світу. У 2005 році Авторська гільдія та кілька окремих авторів подали до суду на Google за порушення авторських прав, стверджуючи, що Google робить несанкціоновані копії їх творів. Google стверджував, що їх проєкт підпадає під виняток "справедливого використання" (fair use), оскільки вони лише показують невеликі уривки тексту у результатах пошуку. Для підтримки своєї думки Google використовував семантичні моделі для аналізу тексту книг та визначення найбільш відповідних частин запитам користувачів. У 2015 році федеральний суд вирішив на користь Google, заявивши, що їх проєкт становить справедливе використання.

У 2009 році Фредрік Кольтінг написав книгу "60 Years Later: Coming Through the Rye", в якій з'явилась персонаж, що базується на Холдені Колфілді з роману Дж.Д. Селінджера "The Catcher in the Rye". Агенти письменника подали позов проти Кольтінга за порушення авторських прав, стверджуючи, що він створив несанкціоноване продовження роману. Кольтінг стверджував, що його книга є твором літературної критики та пародією, які є винятками до закону про авторські права. Для підтримки своєї аргументації Кольтінг використовував семантичні моделі для аналізу тексту обох романів та показу відмінності між ними. У 2010 році федеральний суд виніс рішення на користь майна Селінджера та видав наказ про заборону публікації книги Кольтінга.

Актуальність теми. Використання семантичних моделей має свої обмеження, і їх інструментарій потребує поєднання з іншими методами аналізу для забезпечення більш об'єктивного розуміння тексту. Застосування семантичних моделей у галузі експертизи авторського права та літературного аналізу має великий потенціал для майбутнього, зокрема, для програм виявлення плагіату та цифрових інструментів для допомоги фахівцям з авторських прав. Втім, ефективність семантичних моделей для цих галузей залежить від фахових компетентностей експертів і науковців, як і будь-якої нової технології.

Мета дослідження: оглядово представити провідні семантичні моделі, що застосовуються у літературному аналізі та експертизі об'єктів авторського права, проаналізувати їхню функціональну сутність та специфіку застосування.

Аналіз останніх досліджень і публікацій. Використання семантичних моделей для літературного аналізу та експертизи авторських прав має досить недовгу історію, яка починається з появи так званих «цифрових гуманітарних наук» у 1990-х роках. Одним з перших прикладів використання семантичних моделей для літературного аналізу стала ініціатива ТЕІ (Text Encoding Initiative), започаткована у 1987 році з метою розробки стандарту для кодування літературних текстів у електронному вигляді. ТЕІ використовує мовні маркери для створення структурованих форматів текстів, які можуть бути знайдені, проаналізовані та поширені на різних платформах. У 2000-х роках з'явилося чимало нових проєктів, які використовували семантичні моделі для аналізу літературних текстів. Один з найбільш відомих серед них – Stanford Literary Lab, заснований у 2010 році, який використовував обчислювальні інструменти для аналізу великих наборів даних літературних текстів. Також на початку 2000-х років, галузь експертизи авторських прав почала проявляти інтерес до семантичних моделей як засобу кращого розуміння юридичних наслідків цифрового тексту. Один з перших прикладів цього – проєкт Creative Commons, який розробив набір ліцензій, які можуть бути застосовані до цифрового контенту, щоб вказати, як його можна поділитися та використовувати повторно. Сьогодні семантичні моделі широко використовуються як у літературному аналізі, так і в експертизі авторських прав.

Виклад основного матеріалу дослідження. Розглянемо деякі з магістральних напрямків застосування семантичних моделей у літературному аналізі та експертизах текстів як об'єктів авторського права.

Одна з популярних семантичних моделей, яка використовується у експертизі авторського права, – це латентний семантичний аналіз (LSA). Його сутність полягає у аналізі взаємозв'язків між словами та виявленні їхньої імпліцитної семантики. Ця модель може бути використана для порівняння тексту двох документів та виявлення подібності змісту, навіть якщо використані відмінні слова або фрази. Ідентифікуючи ці схожості, експерти із авторського права можуть визначити, чи є один документ несанкціонованою копією іншого.

Для прикладу, з метою дослідження певного документа, що містить ознаки академічного плагіату, ми можемо використовувати семантичні моделі для аналізу цього документа та порівняння його з відомими джерелами інформації, щоб визначити, чи плагіат дійсно може бути підтверджений. З цією метою LSA застосовується для генерації так званої «векторно-просторової моделі» для документа, яка представляє досліджуваний об'єкт у вигляді набору векторів у багатовимірному просторі, де кожен вектор представляє різні слова або фрази у тексті. Наступним кроком є порівняння векторно-просторової моделі із аналогічними моделями інших відомих документів, таких як раніше опубліковані роботи або інші джерела інформації, щоб виявити ознаки схожості (Manning, Schütze, 1999, p.123).

Подібним чином латентний семантичний аналіз у літературознавстві може бути використаний для визначення теми, аналізу розвитку персонажів та розуміння багаторівневого змісту тексту. Ця модель допомагає літературним критикам і науковцям відшукати такі приховані елементи тексту, які залишаються непоміченими під час першого «поверхневого читання» (що цілком відповідає основоположним принципам «close reading» Нової Критики).

Окрім латентного семантичного аналізу, у літературознавстві використовується так зване тематичне моделювання (topic modeling). Тематичне моделювання працює шляхом аналізу слів, котрі використовуються в тексті та групує їх за темами, заснованими на їх латентному змісті. Ця модель може допомогти літературним аналітикам визначити повторювані теми в тексті, а також простежити взаємозв'язки між різними персонажами і сюжетними точками. Втім, цей науково-методичний засіб використовується й у авторському праві. У галузі експертизи таких об'єктів авторського права як літературні тексти, згадані семантичні моделі можуть бути використані для аналізу значних обсягів тексту та виявлення можливих порушень авторських прав. Ці моделі здатні виявляти схожості між різними текстами, навіть якщо використані слова різні, аналізуючи латентний зміст цих слів (Aggarwal, Zhai, 2012, 2014, p. 208).

Крім LSA й тематичного моделювання, існують й інші семантичні моделі, котрі набули розповсюдження у науковій практиці фахівців експертизи авторського права та літературного аналізу.

Ще однією варіацією семантичних моделей є так звана «класифікація текстів»: техніка, що полягає в призначенні передбачених категорій або міток для документа або набору документів. У контексті літературознавства класифікація текстів може бути використана для категоризації творів літератури за жанром, періодом часу або іншими критеріями. У галузі експертизи авторських прав класифікація текстів може бути використана для категоризації текстів за рівнем їх схожості або для ідентифікації текстів, які можуть порушувати авторські права.

Аналіз настроїв (sentiment analysis) – це специфічний різновид семантичної моделі, яка акцентує емоційний тон певного тексту. У випадку експертизи авторського права, аналіз настроїв може бути використаний для визначення того, чи є документ пародією, або прямим копіюванням вже існуючого твору (Jurafsky, Martin, 2014, p.206-209).

Одна з найбільш поширених моделей штучних нейронних мереж для обробки природної мови – модель Transformer, яка була представлена у статті А.Васвані 2017 року "Увага – все, що вам потрібно", і вона представила архітектуру Transformer для задач моделювання типу sequence-to-sequence в обробці природної мови. Transformer є архітектурою нейронної мережі, яка використовує механізми уваги для обробки вхідних послідовностей та генерації вихідних послідовностей. У статті запропоновано новий підхід до машинного перекладу, який досяг стану мистецтва на кількох тестових датасетах, що при цьому є більш обчислювально ефективним, ніж попередні методи, які ґрунтувалися на рекурентних нейронних мережах. Transformer з тих пір став широко використовуваним моделлю в NLP і був адаптований для різних інших задач, таких як аналіз тональності, мовне моделювання та відповіді на запитання. Ця модель використовує механізми самоуваги для вивчення взаємозв'язків між різними словами в реченні і продемонструвала кращу ефективність, порівняно зі старішими моделями, в різних задачах обробки природної мови.

Хоча семантичні моделі довели свою ефективність як інструменти у сферах експертизи авторських прав та літературного аналізу, втім, вони мають кілька помітних обмежень та потенційних методичних недоліків. Назвемо деякі із основних проблем, що виникають при використанні семантичних моделей.

Ось декілька способів вирішення потенційних проблем з використанням семантичних моделей у галузі авторського права та літературному аналізу:

– комбінування кількісного та якісного аналізу: Щоб уникнути надмірного впливу автоматизованого кількісного аналізу, важливо також включати в процес якісний аналіз, який здійснює людина-експерт. Це передбачає уважне прочитання текстів, врахування ширшого культурного та історичного контексту, фінальну експертну оцінку;

– використання відразу декількох моделей: замість того, що застосовувати лише одну семантичну модель, необхідно використовувати відразу декілька моделей для аналізу певного тексту. Це допоможе подолати обмеження кожної моделі та забезпечити більш детальний та комплексний підхід;

– вирішення проблеми надмірного впливу окремого комплексу даних: важливо переконатися, що набір даних, який використовується для навчання семантичної комп'ютерної моделі, є різноманітним та репрезентативним. Це передбачає збирання даних із різних джерел та використання варіативних методів для забезпечення рівноваги репрезентації різних груп у наборі даних алгоритму комп'ютерної моделі.

– студіювання ширшого контексту: хоча семантичні моделі здатні ефективно ідентифікувати подібні фрагменти в тексті, вони не завжди можуть врахувати культурний або історичний контекст, в якому розташований досліджуваний текст в цілому. Щоб вирішити цю проблему, важливо враховувати контекст, в якому було створено текст, та обов'язково залучати до формування кінцевих висновків проміжні умовиводи людини-експерта.

– верифікація результатів експертизи комісією спеціалістів: для уникнення хибної інтерпретації, результати семантичного аналізу можуть бути розглянуті групою спеціалістів, котрі добре обізнані із контекстом досліджуваних об'єктів. Це може допомогти забезпечити особливу точність аналізу та виявити важливі нюанси розуміння тексту.

Висновки. Представлений у публікації огляд сутності та особливостей використання семантичних моделей у експертизі авторського права та літературному аналізу дозволяє зробити наступні висновки. Семантичні моделі значно поліпшують ефективність аналізу текстів. Ці моделі використовують передові технології для ідентифікації структурних патернів, тем і схожих елементів у текстах, а також дозволяють фахівцям отримати глибше розуміння творів, які вони аналізують.

Хоча семантичні моделі мають багато переваг, вони також мають свої функціональні обмеження. Важливо бути обізнаним із цими обмеженнями та використовувати аналітичні моделі у поєднанні з іншими методами для забезпечення більш комплексного розуміння тексту.

Зрештою, використання семантичних моделей в експертизі авторського права та літературному аналізу має величезну практичну цінність, але важливо бути обізнаним з їх обмеженнями та використовувати їх разом з іншими методами аналізу для забезпечення максимального об'єктивного і всебічного дослідження. Із розвитком технологій будуть розроблені нові й досконаліші семантичні моделі, які дозволять проводити ще більш детальний аналіз текстів.

Література:

1. Landauer, T. K., Foltz, P. W., & Laham, D. (1998). An introduction to latent semantic analysis. *Discourse Processes*, 25(2-3), 259-284.
2. Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3(Jan), 993-1022.
3. Manning, C. D., & Schütze, H. (1999). *Foundations of statistical natural language processing*. MIT Press.
4. Aggarwal, C. C., & Zhai, C. (2012). *Mining text data*. Springer Science & Business Media.
5. Jurafsky, D., & Martin, J. H. (2014). *Speech and language processing: An introduction to natural language processing, computational linguistics, and speech recognition*. Pearson Education.
6. Rosenberg, M. (2015). *Text mining with R: A tidy approach*. O'Reilly Media, Inc.
7. McMenamin, G. R. (2017). *Basic text mining: classify, cluster, visualize data*. CRC Press.