

Вып. 4 (93). – С. 3–16. 6. Коряшкіна Л. С., Михалева А. А., Навоенко В. И. Применение методов оптимального разбиения множеств к непрерывным задачам многократного покрытия // *Питання прикладної математики і математичного моделювання: зб. наук. праць.* – Дніпропетровськ, 2014. – С. 141–154. 7. Киселева Е. М., Коряшкіна Л. С., Михалева А. А. Непрерывная задача многократного шарового покрытия с ограничениями и метод ее решения // *Системні технології. Дніпропетровськ.* – 2015. – №1. – С. 165–179. 8. Z. Drezner. The p -centre problem – heuristic and optimal algorithms. *J/ OR Soc.* 1984. V. 35. P. 741 – 748. 9. Галиев Ш. И. Направление убывания для минимаксиминных задач // *Журн. вычисл. матем. и матем. физ.* – 1994. – Т. 34. – № 3. – С. 323–343. 10. Preparata F. P., Shamos M. I. 1985. *Computational Geometry: An Introduction (Texts and Monographs in Computer Science)*. New York: Springer-Verlag New York, Inc: 390. 11. Коряшкіна Л. С. Обобщение одного класса задач бесконечномерного математического программирования // *Математичне та імітаційне моделювання систем. МОДС 2015: тези доповідей Десятої міжнар. наук.-практ. конф. (Чернігів, 22 – 26 червня 2015 р.).* – Чернігів: ЧНТУ, 2015.– С. 160–164. 12. Шор Н.З. *Методы минимизации недифференцируемых функций и их приложения.* – К.: *Наук. думка*, 1979. – 200 с.

УДК 519.8

О. Мриглод

Інститут фізики конденсованих систем НАН України

АВТОМАТИЗОВАНИЙ АЛГОРИТМ ПОШУКУ ТЕРМІНІВ У НАУКОВИХ ПУБЛІКАЦІЯХ

© Мриглод О., 2015

Описано послідовність застосування одного з алгоритмів автоматизованого пошуку наукових термінів, модифікованого з огляду на специфіку поставленої задачі. Проаналізовано сукупність наукових документів з вибраної тематики, погрупованих за кількома дисциплінарними напрямками. В результаті комбінації лінгвістичного та статистичного підходів до аналізу текстів визначено перелік найважливіших термінів, що дають змогу оцінити спектр дрібніших тематичних напрямів у публікаціях з кожної дисципліни.

Ключові слова: інтелектуальний аналіз тексту, автоматизований пошук термінів, текст, публікація.

The application of partially modified semi-automatic algorithm of scientific terms searching is described in this paper. The set of research papers of a given topic within several disciplines were analyzed. The combination of linguistic and statistical approach to the analysis of texts gave a possibility to get the list of the most important terms. These terms can be used to reveal the spectra of subtopics in the set of selected publications within each discipline.

Key words: text mining, semi-automatic terms identification, text, publication.

Вступ

Серед наукометричних досліджень важливе місце посідають проблеми вивчення структури науки та її еволюції. Виявлення так званих “гарячих напрямів” та спостереження за розвитком окремих тематик – це задачі, розв’язок яких може бути вельми корисним для практичного використання. Адже інформація про те, які напрями у науці сьогодні є особливо затребуваними та

гіпотетично перспективними, є необхідною для прийняття рішень, починаючи від постановки задачі для майбутніх дослідників і завершуючи розподіленням державних або грантових коштів. Завдання моніторингу наукових напрямів, на перший погляд, може видатись достатньо простим, проте з огляду на складність системи науки та процесів, що відбуваються у ній, однозначного вирішення досі не має [1, 2]. Для того, щоб погрупувати публікації або, скажімо, видання за тематичною ознакою, і тим самим визначити певну структуру наукових досліджень, використовують різні методи. Окрім експертного аналізу, тобто класифікації чи сортування “вручну”, використовуються алгоритми автоматизованої кластеризації на основі даних про співавторство або ж співцитування (про одні із перших спроб див. у [3, 4]). Також проблема визначення тематичного забарвлення наукових текстів суміжна з проблемою виділення основних тематичних концепцій та побудови тезаурусів, що часто використовує лінгвістичні підходи для аналізу власне змісту текстів – так званого контент-аналізу. А вже підзадача виявлення тематики публікації виявляється потенційно корисною для цілого спектра практичних застосувань: організації релевантного інформаційного пошуку, каталогізації та рубрикації електронних ресурсів або ж публікацій у виданні, автоматичного пошуку рецензентів та багатьох інших. Загалом, завдання зводиться до розроблення методів автоматичного чи хоча б автоматизованого (за часткової участі людини-експерта) аналізу текстів з виділенням ключових тематичних концепцій, поданих у вигляді ключових слів. У випадку наукових текстів результатом може бути перелік значущих наукових термінів (формальніше визначення *терміна* розглянуто далі), що, власне, відображають ці концепції.

Коротко про підходи до виявлення ключових слів у текстах

Розробити автоматичний чи хоча б напівавтоматичний спосіб виділення ключових термінів, що описували б основні концепції документів, намагаються вже не перше десятиліття (див., наприклад, [5–8]). Адже така задача потенційно має не одне практичне застосування: тематичне маркування виявлених груп документів, завдання інформаційного пошуку, каталогізації, та інші, згадані вище.

Перш ніж перейти власне до обговорення способів аналізу текстів, необхідно визначити, що ж ми розуміємо під *терміном*. Насправді не існує чітко формалізованого визначення, проте зазвичай термінами називають так звані “змістовні” (чи “сигнальні” [5]) слова або словосполучення, що передають основні змістові ідеї тексту, тобто відображають ту чи іншу тематичну концепцію, висвітлену в документі. На відміну від таких “змістовних” слів, “функціональні” використовуються для зв’язування речень та передавання додаткової інформації [9]. Можна уявити, що “функціональні” слова є середовищем, яке забезпечує розташування “змістовних” слів – наче риб у воді. До “функціональних”, зокрема, зараховують усі види сполучників та службових слів.

Важливо розуміти, що “змістовність” кожного конкретного слова невід’ємна від контексту. Так, слово “шум” може бути вторинним у тексті з біології, проте стати терміном для фізичної публікації. Релевантність слова може змінюватись навіть у межах однієї дисципліни, тому множина термінів завжди є індивідуальною для конкретно вибраного набору документів.

Окремою проблемою є охоплення аналізом не лише одиничних, але й складених термінів (з кількох слів). Їх автоматичне виділення технічно є проблематичнішим, проте часто вони допомагають уточнити загальніші за значенням одиничні терміни, детальніше описати ту чи іншу концепцію документа. Скажімо, коли іменник “пухлина” може означати доволі широкий медичний спектр тематик, то словосполучення “тироїдна пухлина” вже значно звужує коло пошуку. У цьому випадку знову не знімається питання про те, що в подальшій роботі все-таки вважати терміном: слово “пухлина” чи словосполучення “тироїдна пухлина”. Тут, як правило, потрібно приймати рішення знову ж таки для кожного конкретного набору документів.

Вважається, що найнадійнішим методом визначення множини ключових термінів для корпусу документів чи певної галузі є залучення експертів – фахівців у відповідній ділянці. Автоматизувати цей процес поки що не вдається власне через відсутність абсолютних критеріїв, багатозначність

мови, її контекстність тощо. Проте вже тривалий час пропонуються та вивчаються напівавтоматичні методи аналізу текстів та визначення ключових слів, термінів або ж концептів. Для знаходження множини слів чи словосполучень, що потенційно можуть бути такими ключовими словами, можна використовувати різні принципові підходи: на основі статистичного, синтаксичного чи змішаного аналізу слів у наборі документів [6, 9]. У першому випадку текст розглядається лише як випадковий набір або ж впорядкована послідовність елементів – слів. Тоді можна робити частотний аналіз, знаходити типові послідовності елементів та застосовувати інші статистичні підходи. У другому випадку враховуються синтаксис, частини мови та структура слів у реченнях. А, зрештою, на практиці найчастіше використовується комбінація цих двох методів [9, 11].

Якщо знехтувати структурою документа та вважати його “мішком зі словами”, то можна проаналізувати частоту вживання слів k , побудувавши її розподіл. Ще у середині минулого століття доведено, що в результаті такого аналізу одержимо степеневий закон розподілу слів, відомий як закон Зіпфа [10]. Останній говорить про те, що у тексті типово є велика кількість різних слів, що трапляються один раз або кілька разів, і лише декілька таких, що вживаються дуже часто (див. далі рис. 3). Певний парадокс полягає у тому, що хоч частота вживання слів є одним із базових понять, проте на основі лише частотного розподілу неможливо визначити, які ж зі слів можна вважати ключовими або такими, сукупність яких описує власне основні концепції документа. Найчастіше вживані слова, як правило, є дуже загальними за змістом, а рідко вживані – дуже конкретними, проте не можуть вважатися статистично значущими. Вважають, що найкращі кандидати у “змістовні” слова розмістяться десь посередині згаданого частотного розподілу (див. рис. 3) – такі, що вживаються не найчастіше, проте і не надто рідко [5, 9, 11]. Залежно від зростання довжини текстів частота “функціональних” та “змістовних” слів змінюється по-різному: для перших вона пропорційно зростає, тоді як для других просто відбувається розширення словника (більший текст – більша імовірність обговорення нової концепції – нові ключові слова/терміни) [9, 11].

Окрім частоти вживання, додаткову інформацію можна отримати, враховуючи структуру корпусу та окремих його документів. Відомо, що «функціональні» та «змістовні» слова неоднаково розподілені серед документів чи структурних частин одного документа, тоді як перші характеризуються швидше рівномірною розкиданістю по корпусу та документах, другі вживаються нерівномірно – сконцентровано у групі документів (чи в певному місці окремого документа) та рідко в інших.

Вже побіжний огляд показує велику кількість неоднозначностей, пов’язаних із намаганням автоматично виявити значущі слова, що б відображали тематичний спектр набору документів. При цьому досі йшлося лише про формування списку кандидатів на терміни – слів/словосполучень, які надалі потрібно оцінити на предмет їх “значущості”, тобто релевантності, специфічності та важливості для конкретного набору текстів. Знову ж таки, за відсутності визначення такої значущості, можна використовувати різні методи для “зважування” потенційних термінів і порівняння їх ваг між собою. Наприклад, можна використати вже згадану властивість неоднорідного розподілу “змістовних” слів. Про міру “специфічності” слова/словосполучення може говорити добуток частоти його вживання k_i на так звану обернену частоту документів idf_i , в яких воно трапляється. Остання величина дорівнює відношенню загальної кількості документів у корпусі N до кількості документів, в яких трапилось це слово/словосполучення n_i : з огляду на типово велике значення N , можна розглядати логарифм відношення, тобто $idf_i = \log(N/n_i)$. З одного боку, величина $k_i \times idf_i$ буде пропорційною до загальної вживаності кандидата у терміни, а з іншого боку – обернено пропорційною до кількості різних документів, де він вживається [9]. Існує ціла низка методів зважування та нормалізації, що можуть враховувати довжини документів чи їх

структуру (тобто місця локалізації слів у певних частинах документа), факти співпояв слів, їх контекст тощо.

Якщо йдеться виключно про наукові статті та виявлення у них наукових термінів, то задача дещо полегшується з огляду на чітку структурованість таких документів. Вважається, що значущі терміни найбільше сконцентровані у певних структурних частинах, таких як заголовки, анотація, перший абзац статті чи висновки. Деякі дослідники вважають, що комбінація заголовка та анотації є достатньою основою для аналізу (наприклад, [6, 12]), інші вважають найрелевантнішими окремо взяті заголовки (див. [13]), проте завжди потрібно враховувати невеликий розмір цих фрагментів з погляду статистичного підрахунку частот. Крім того, значущість тих чи інших структурних елементів статей є різною для природничих та гуманітарних наук.

Можна виділити типові риси наукових термінів загалом. Так, вважається, що найчастіше це одноосібні іменники або ж словосполучення, що формуються навколо головного іменника, доповнюючись іншими іменниками, прикметниками тощо. Потрібно зауважити, що такі висновки поки що найобгрунтованіші для англійських текстів як найдослідженіших [9, 11].

Далі у роботі описано покрокову напівавтоматичну процедуру виявлення наукових термінів у публікаціях вибраної тематики. Таке завдання поставлено в межах ширшої задачі дослідження реакції наукової спільноти – що відображається власне в опублікованих роботах – на визначену подію.

Постановка задачі: приклад тематичного аналізу наукових публікацій

Нещодавно під час вивчення реакції наукової спільноти на Чорнобильську аварію [14, 15] зібрано бібліографічні дані про усі релевантні до проблеми наукові публікації в базі даних Scopus (www.scopus.com) на початок 2015 р. Загалом кінцевий перелік містив понад 9,5 тис. бібліографічних записів про публікації, що містили різні написання слова “Chornobyl” у заголовках, анотаціях чи ключових словах. Зібрано тематичну колекцію наукових документів із вузької тематики. Окрім дослідження розподілу публікацій за галузями науки та за роками, аналізу відповідної мережі співпраці на рівні країн, виявлення зміни зацікавленості в межах різних дисциплін та інших завдань, цікаво було дослідити тематичний спектр всередині зібраного корпусу документів. Адже, поряд із загальнодисциплінарними тенденціями до підвищення чи загасання інтересу в межах тієї чи іншої галузі науки, можна очікувати зміни тематичного спектра на тоншому масштабі, в межах однієї дисципліни – адже з часом актуальність одних проблем втрачається, тоді як інші починають активно досліджуватись. Такий детальніший аналіз тим цікавіший з огляду на те, що сьогодні для домінуючих дисциплін (за кількістю чорнобильських публікацій у Scopus) спостерігається більш-менш стала картина, тобто щорічна кількість публікацій коливається навколо певного значення. З іншого боку, для низки інших дисциплін спостережено тенденції до загасання (скажімо, для ветеринарії) чи зростання (наприклад, для економіки та фінансів) інтересу до чорнобильської тематики [14, 15]. Поставлене завдання виявлення внутрішніх тематик для набору статей хоча б для п'яти найпоширеніших у базі Scopus дисциплін: медицини (3 635 статей); наук про навколишнє середовище (3 156); енергетики (1 470); фізики та астрономії (1 437); біохімії, генетики та молекулярної біології (1 198).

Виявлення ключових термінів дає змогу визначити найактуальніші завдання за часом. Такий аналіз також допомагає побудувати карту наукових дисциплін – тобто візуально показати, як знайдені терміни (а отже, і тематичні піднапрями) взаємопов'язані між собою. Наприклад, на рис. 1 показано карту термінів для нашого набору статей з дисциплін, які домінують, згенеровану за допомогою спеціальної програми VOSviewer [16, 17]. У цьому випадку для автоматичного виділення термінів використовувалась інформація про їх співпояви у заголовку та анотації (що в цьому випадку трактується як один документ) кожної статті. Вбудований алгоритм кластерування дає змогу розрізнити декілька тематичних груп: чотири більші та дві менші. І хоча ці групи не відповідають взаємно однозначно п'яти дисциплінам, що домінують, можна чітко розрізнити дві найбільші ділянки досліджень, що стосуються аварії на Чорнобильській АЕС : вплив на здоров'я

людини (три великі кластери зліва) та наслідки для навколишнього середовища (великий кластер справа). Також достатньо добре візуально розрізняються піднапрями, що пов'язані із онкологічними захворюваннями, генетичними ефектами та аналізом шляхів забруднень у різних середовищах.

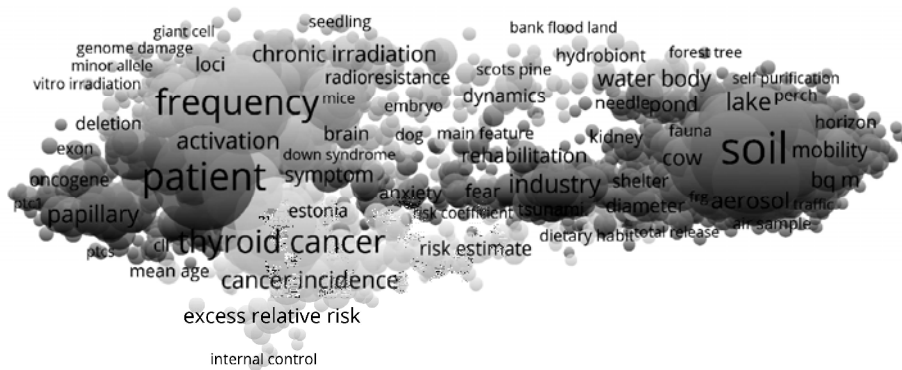


Рис. 1. Карта термінів для чорнобильських публікацій з медицини; наук про навколишнє середовище; енергетики; фізики та астрономії; біохімії, генетики та молекулярної біології, знайдених у базі Scopus на початок 2015 р. Показано 60 % найбільш релевантних термінів

Покрокова процедура виявлення термінів

Для того, щоб отримати ключові слова, набір яких найточніше описує загальний тематичний зміст чорнобильських наукових публікацій з п'яти дисциплін, що домінують, перелічених вище та двох додаткових – мистецтва (лише 59 публікацій у Scopus) та соціології (310) як яскравих представників гуманітарних наук, – модифіковано процедуру, запропоновану в [11]. Основні кроки описано нижче.

Крок 1. Передусім потрібно визначити, що буде основою для аналізу: окрім інформації про авторів та видання, у нашій базі зібрані анотації, заголовки та ключові слова. Оскільки ключові слова за замовчуванням можна вважати авторськими термінами, за основу об'єктивного аналізу беремо анотації та/чи заголовки. Заголовок за призначенням мав би найточніше та найкоротше вказувати на суть статті, проте все частіше заголовки покликані скоріше привабити читача, аніж вказати на зміст. Тому не комбінувалися заголовок+анотація, а порівнювалися два окремі випадки, коли як *документи* трактували окремо заголовки та окремо анотації.

Крок 2. Для того, щоб одержати інформацію про лінгвістичні властивості слів у документах, можна використати одну із доступних у вільному доступі програм – а саме TreeTagger [18]. Ця програма призначена для маркування слів у тексті за частинами мови та знаходження базової форми для кожного слова – що пізніше дає змогу нехтувати різними закінченнями, що залежать від числа або роду. Приклад результатів обробки фрагмента тексту програмою наведено на рис. 2.

```
colony-stimulating|NN|colony-stimulating
factors|NNS|factor
for|IN|for
the|DT|the
treatment|NN|treatment
of|IN|of
the|DT|the
hematopoietic|JJ|hematopoietic
component|NN|component
```

Рис. 2. Результат обробки програмою TreeTagger фрагмента тексту: “colony-stimulating factors for the treatment of the hematopoietic component”

Крок 3. Із одержаного переліку легко відібрати потрібні конструкції – які ми далі називатимемо *семантичними одиницями*, дотримуючись означення, введеного в роботах [11, 15, 19] (semantic units) – одиничні іменники або ж словосполучення, що складаються лише з іменників та прикметників. Подібно, як у [11, 19], використано загальне правило для відбору складених конструкцій: **прикметник *іменник*. Ця загальна форма означає, що фраза може розпочинатись із довільної кількості прикметників та завершуватись довільною кількістю іменників, наприклад: “post-chernobyl radioactive contamination”, “radionuclide contamination source”. Крім того, певні найпростіші правила перетворення дають змогу не відсіювати ті конструкції, які існують неявно:

- Ї Фраза, що містить на початку декілька прикметників, розділених комами, та один іменник в кінці, перетворюється на декілька фраз, що складаються з одного прикметника та іменника; наприклад, конструкція “personal, political, linguistic, historical complexity” буде врахована у вигляді чотирьох семантичних одиниць: “personal complexity”, “political complexity”, “linguistic complexity”, “historical complexity”.
- Ї Фраза, у якій один прикметник відноситься до декількох іменників, розділених сполучником “and” (“і” чи “та” з англ.), перетворюється на декілька фраз, що містять прикметник та один із іменників; наприклад, із конструкції “individual responses and responsibilities” отримаємо “individual responses” та “individual responsibilities”.

Звичайно, неможливо передбачити та забезпечити перетворення усіх можливих варіантів конструкцій коректно. Поряд із технічними помилками чи описками існують такі, які за змістом мали б розглядатись, проте не відповідають заданому шаблону, містячи сполучники. Скажімо, класичним прикладом є словосполучення “degrees of freedom”. Тому потрібно пам'ятати про неминучі похибки автоматичних процедур, особливо тих, що стосуються обробки природної мови.

Відбираючи семантичні одиниці, відразу ж варто відсіяти такі, що занесені у так званий стоп-список. Цей перелік, як правило, складається вручну і конкретно для кожного набору документів. Так, наприклад, у нашому випадку стоп-словами стали семантичні одиниці “article” (“стаття”), “Elsevier” (назва видавництва), різноманітні одиниці вимірювання величин тощо.

У результаті виконання перших трьох кроків для нашого корпусу отримано список із 80 094 семантичних одиниць для анотацій та 15 020 – для заголовків.

Крок 4. Наступним кроком є підрахунок семантичних одиниць для проведення частотного аналізу. Проте проста, на перший погляд, процедура не є однозначною.

- Ї Передусім, всупереч поширеній практиці, що передбачає відсіювання маловживаних семантичних одиниць одразу ж [11, 19], ми їх спочатку не вилучаємо. Оскільки розміри наших вибірок даних, особливо для гуманітарних дисциплін, є невеликими, таке відсіювання завадить побачити загальну форму частотного розподілу. На рис. 3 зображено частотно-рангові розподіли семантичних одиниць: для їх побудови спочатку треба посортувати усі семантичні одиниці за зменшенням частоти появ k (по вертикальній осі на рисунку), а тоді присвоїти їм відповідні ранги r у послідовності зростання (горизонтальна вісь, відповідно). Отже, на рис. 3 бачимо відповідні графіки для анотацій (а) та заголовків (б), що демонструють близькість до степеневого закону. Така форма розподілу типова для текстів, написаних природною мовою [10, 20]. І навіть більше, не лише форма, але й нахил одержаної кривої (зі значенням експоненти, що близька до -1) свідчить про універсальні характеристики таких коротких текстів, як

анотації (рис. 3, а). Такий результат не є цілком очевидним, адже в цьому випадку ми оперуємо не частотою слів, а частотою лише певних семантичних конструкцій.

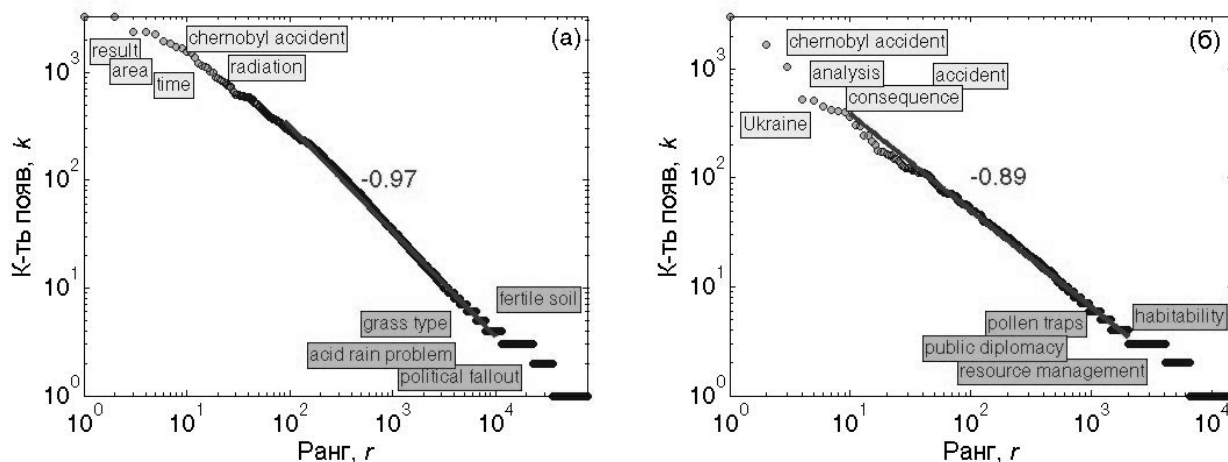


Рис. 3. Частотно-рангові розподіли семантичних одиниць для (а) анотацій та (б) заголовків чорнобильських публікацій з семи вибраних дисциплін у базі даних Scopus. Розподіли добре апроксимуються степеневими залежностями $k \sim r^{-a}$ із показниками $a \approx 0.97$ та $a \approx 0.89$, відповідно

Ї Перед тим, як рахувати частоту появ складених конструкцій, тобто тих, що складаються з двох чи більше слів, можна спочатку оцінити їх зв'язаність, статистично підтвердивши, що ціла конструкція трапляється достатню кількість разів порівняно з частотою її окремих елементів. Тобто можна йти шляхом первинної перевірки того, чи знайдені складені семантичні одиниці повинні фігурувати як одне ціле, чи доречно їх розділити на менші елементи (див. [11, 19]). Проте у нашому дослідженні такої перевірки не зроблено, оскільки розраховується частота не лише всієї складеної конструкції, але й її складових: окремо останнього слова (іменник, що трактується як головний у семантичній одиниці), а далі останнього слова разом із тими, що йому передують, додаючи їх один за одним. Наприклад, для семантичної одиниці “low-dose radiation exposure” підраховується частота вживань “exposure”, “radiation exposure” та “low-dose radiation exposure”. Тобто ми беремо до уваги так звані вкладені терміни (nested terms): терміни, що можуть бути самостійними або входити до складу інших [11].

Ї Для підрахунку кількості вживань семантичної одиниці використано так званий бінарний спосіб, тобто якщо в тому самому документі вона виявлена більше від одного разу, все одно “зараховується” лише одноразово. Отже, загальна кількість появ дорівнюватиме кількості документів, у яких трапилася ця конструкція.

Крок 5. Нарешті, із загального списку зібраних та підрахованих семантичних одиниць потрібно виділити власне *терміни*. Як вже згадано вище, важливість семантичної одиниці як терміна не є прямо пропорційною до її частоти вживання. Найчастіше вживані слова, що потрапляють на початок частотно-рангового розподілу типу Зіпфа, є, як правило, найзагальнішими за змістом для цього набору документів. Скажімо, у нашому випадку це такі семантичні одиниці, як “Chornobyl accident”, “radiation”, “Ukraine” тощо (рис. 3). З іншого боку, слова із найменшою кількістю вживань, що потрапляють у “хвіст” розподілу, насправді вельми специфічні, проте не розглядаються як статистично значущі: наприклад, “grass type”, “public diplomacy” тощо. Найвірогіднішими кандидатами у терміни є семантичні одиниці, розміщені посередині розподілу (рис. 3) [9]. Тому на цьому кроці важливо застосувати певні частотні (чи інші) фільтри для

відокремлення термінів. Спочатку введемо нижнє критичне значення частоти $k_c = 4$, тобто відсіємо усі семантичні одиниці, що виявлені у менше ніж чотирьох документах. Це значення вибрано за допомогою емпіричного спостереження: частота появ семантичних одиниць, якщо $k \leq k_c$ знижується повільно, тоді як вище від цього значення режим змінюється на значно швидший.

Черговий фільтр пов'язаний із нерівномірним поширенням термінів між дисциплінами. Тоді як семантичні одиниці можуть однаково активно використовуватися у статтях, що стосуються різних галузей знань, термінами називатимуться ті, що характерні для певної дисципліни чи кількох дисциплін. Щоб виразити концепцію такого нерівномірного розподілу в числовому вигляді, використаємо ідею так званої “термінності” (*termhood*), запропоновану у [11, 19]. Йдеться про міру специфічності семантичної одиниці, тобто її характерності для певної/певних дисциплін. Для розрахунку термінності спочатку потрібно побудувати ймовірнісний розподіл частоти вживання усіх кандидатів у терміни $P(d)$ за дисциплінами ($d = 1..7$). Відповідний графік для бази даних анотацій (майже збігається із аналогічним графіком для заголовків) показано на рис. 4 (лінія з кружечками).

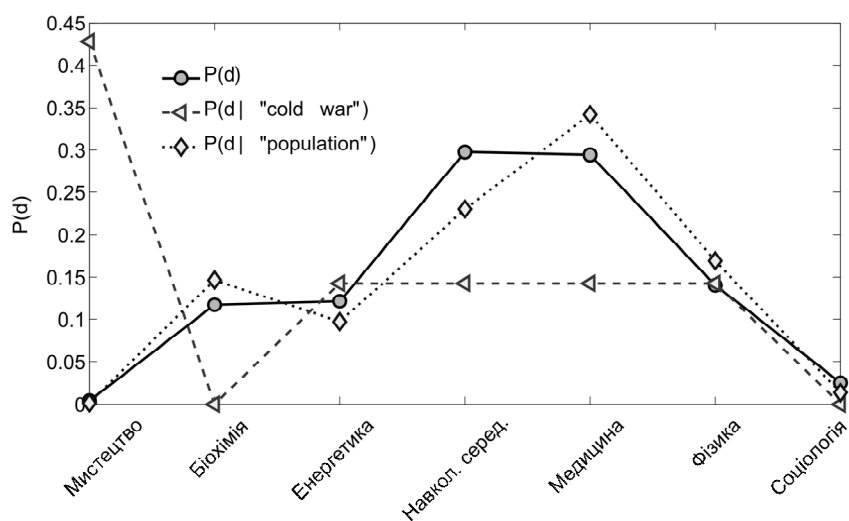


Рис. 4. Ймовірнісні розподіли частоти появ семантичних одиниць у дисциплінах на основі бази даних анотацій чорнобильських публікацій: загальний розподіл $P(d)$ та два індивідуальні розподіли для семантичних одиниць “cold war” та “population”

Із рис. 4 видно, що найбільша кількість семантичних одиниць стосується наук про навколишнє середовище та медицини. Очікувано, що найменші значення відповідають обидвом гуманітарним дисциплінам. Далі аналогічно для кожної семантичної одиниці s (s змінюється від 1 до загальної кількості семантичних одиниць у сформованому списку) будується її власний, індивідуальний розподіл $P(d | s)$. Інакше кажучи, $P(d)$ показуватиме ймовірність будь-якої семантичної одиниці “потрапити” у певну дисципліну d , а $P(d | s)$ – цю ймовірність для конкретної семантичної одиниці s див. рис. 3 (лінії з трикутниками та ромбами). Різниця між загальним $P(d)$ та конкретним $P(d | s)$ розподілом, виражена у числовому вигляді, і буде шуканою величиною. Існують різні математичні способи порівняння розподілів між собою, ми використали запропонований у [11, 19]. Для розрахунку рівня “неподібності” між кожною парою розподілів $P(d)$ та $P(d | s)$ вживається поняття так званої від’ємної ентропії. Так, міру termhood для вибраної семантичної одиниці j розраховують як:

$$t_s = \sum_{d=1}^7 \log p_d, \quad \text{де} \quad p_d = \frac{P(d | s) / P(d)}{\sum_{d'=1}^7 P(d' | s) / P(d')}$$

приймаючи, що $0 \log 0 = 0$. Що вище значення величини t_s , то специфічнішою (характернішим для певної дисципліни чи кількох дисциплін) вважається семантична одиниця – то більше підстав її вважати терміном. На рис. 5 продемонстровано, як для кожної семантичної одиниці змінюються її загальна частота k_s та специфічність t_s .

Очевидно, що необхідно знайти компроміс між цими двома величинами, проте будь-яке рішення буде певною мірою умовним або суб'єктивним. Така ситуація доволі типова для задач, що потребують участі експертів та не можуть бути повністю автоматизованими. У цьому випадку до термінів зараховано ті семантичні одиниці, що відповідали таким критеріям:

- $k_s > k_c$, де $k_c = 4$ (див. вище).
- $t_s > t_c$, де t_c дорівнює медіані.
- Семантична одиниця належить до перших 50 у переліку, відсортованому за значенням величини $t'_s \cdot k'_s$ (добутком t_s та k_s , нормованими так, щоб належати інтервалу $[0..1]$) – така вага дає додаткову перевагу специфічності термінів порівняно з їх частотою.

Деякі семантичні одиниці вилучені зі списку термінів вручну на фінальній стадії – як пояснювалося вище, обійтися без втручання людини-експерта поки що неможливо.

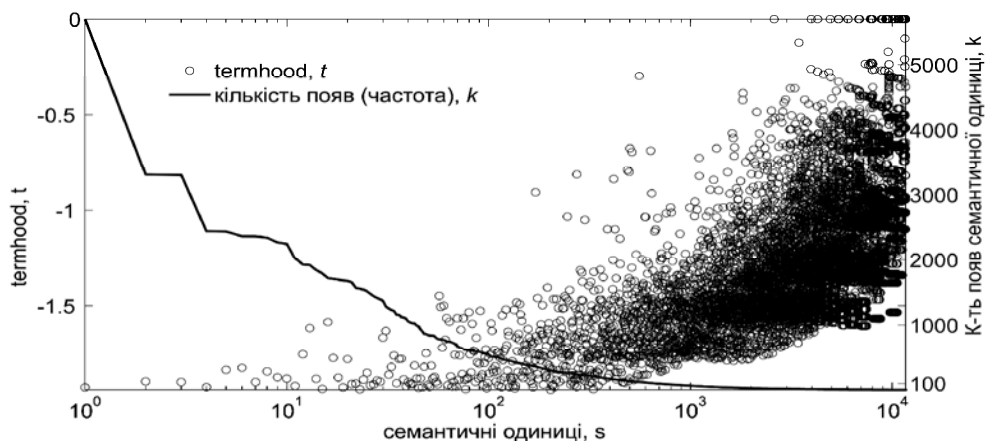


Рис. 5. Специфічність (termhood) t_s відносно частоти появ $k_s > k_c$ для семантичних одиниць на основі бази анотацій

У результаті одержано переліки термінів, що характеризують публікації на тему Чорнобильської аварії для кожної із семи вибраних дисциплін. На їх основі можна робити висновки про піднапрями, актуальні в межах ширших областей досліджень. У табл. 1 та 2 наведено по двадцять найспецифічніших термінів, отриманих на основі анотацій та заголовків, відповідно. Можна побачити, що у часовий період, найближчий до аварії, виділялися терміни із наук про навколишнє середовище. Більшість термінів, характерних для біохімії, генетики та молекулярної біології, вперше з'являються у публікаціях на початку 90-х років. Гуманітарні терміни починають виникати ще пізніше (2002–2006). Це підтверджує думку про те, що чорнобильська тематика досліджувалась в межах різних дисциплін не синхронно [14, 15]. Природно, що безпосередньо після катастрофи акцентували на найшвидших її наслідках для навколишнього середовища, здоров'я людей; з часом все актуальнішими стали віддаленіші наслідки, наприклад, генетичні та онкологічні; натомість після кількох десятиріч обговорюються також економічні, соціальні та культурні проблеми, пов'язані з аварією на ЧАЕС.

**Перша двадцятка термінів, найспецифічніших для чорнобильських публікацій
у межах семи досліджуваних дисциплін, відібраних на основі анотацій статей у Scopus**

Терміни (мовою оригіналу)	Терміни (український переклад)	Характерні для:	Рік першої появи у базі даних
1) carcinoma	1) карцинома	біохімія	1992
2) thyroid carcinoma	2) карцинома щитовидної залози	біохімія	1992
3) tumor	3) пухлина	біохімія	1994
4) gene	4) ген	біохімія	1987
5) rearrangement	5) перебудова	біохімія	1993
6) papillary thyroid carcinoma	6) папілярний рак щитовидної залози	біохімія	1995
7) ptc*	7) ptc	біохімія	1992
8) papillary carcinoma	8) папілярна карцинома	гуманітарні	2002
9) science	9) наука	біохімія	1994
10) carcinogenesis	10) канцерогенез	біохімія	1992
11) malignancy	11) злоякісність	навкол. сер.	1987
12) activity ratio	12) коефіцієнт активності	гуманітарні	2006
13) threat	13) загроза	біохімія	1992
14) metastasis	14) метастази	біохімія	1994
15) surgery	15) хірургічна операція, хірургія	соціологія	1989
16) policy	16) політика	біохімія	1999
17) cleanup worker	17) працівник з очищення (ліквідатор)	біохімія	1997
18) high frequency	18) висока частота	біохімія	1990
19) nuclear disaster	19) ядерна катастрофа	соціологія	1990
20) discharge	20) розряд	навкол. сер.	1988

* ptc – аббревіатура від “papillary thyroid carcinoma”

**Перша двадцятка термінів, найспецифічніших для чорнобильських публікацій
у межах семи досліджуваних дисциплін, відібраних на основі заголовків статей у Scopus**

Терміни (мовою оригіналу)	Терміни (український переклад)	Характерні для:	Рік першої появи у базі даних
1) carcinoma	1) карцинома	біохімія	1993
2) thyroid carcinoma	2) карцинома щитовидної залози	біохімія	1993
3) patient	3) пацієнт	біохімія	1991
4) rearrangement	4) перебудова	біохімія	1991
5) sediment	5) осад	навкол. сер.	1987
6) papillary thyroid carcinoma	6) папілярний рак щитовидної залози	біохімія	1995
7) transport	7) перенесення, транспорт	навкол. сер.	1987
8) mutation	8) мутація	біохімія	1989
9) tumor	9) пухлина	біохімія	1994
10) cleanup worker	10) працівник з очищення (ліквідатор)	біохімія	1993
11) cleanup	11) очищення (ліквідація)	медицина	1992
12) unit	12) модуль	енергетика	1982
13) history	13) історія	гуманітарні	2009
14) pond	14) ставок	навкол. сер.	1987
15) policy	15) політика	соціологія	1988
16) prevalence	16) поширеність, розповсюдження	біохімія	1995
17) Black sea	17) Чорне море	навкол. сер.	1987
18) thyroid disease	18) хвороба щитовидної залози	біохімія	1991
19) radiation protection	19) захист від радіації	гуманітарні	2006
20) forest ecosystem	20) екосистема лісу	навкол. сер.	1991

Висновки

У результаті виконаної роботи можна зробити дві групи висновків. Перша стосується самої процедури виокремлення термінів (ключових/значущих слів) у наукових текстах. Сьогодні є усі підстави вважати, що її повна автоматизація неможлива – на тому чи іншому етапі необхідно залучати експертів у відповідній галузі знань для кожного конкретного набору наукових публікацій. Така експертна участь необхідна як на проміжних стадіях, скажімо, для формування списку слів, які завідомо не є змістовними – стоп-списку, так і на кінцевій стадії для верифікації результатів. Проте численні напрацювання у цьому напрямі забезпечили чималий арсенал підходів та методів для автоматизації окремих етапів процедури пошуку термінів. Звичайно, можна повністю покластися на розроблені програми (такі продукти вже існують, наприклад, VOSviewer [16]), проте необхідно допускати відповідну похибку в результатах.

У роботі реалізовано алгоритм пошуку наукових термінів до бібліометричної бази статей, що стосуються аварії на Чорнобильській АЕС. Цей алгоритм ґрунтується на комбінації лінгвістичних та статистичних методів опрацювання наукових текстів. З врахуванням усіх нюансів та часткою суб'єктивності вдалося сформувати перелік термінів, характерних для п'яти дисциплін із найбільшою кількістю публікацій від 1986 до початку 2015 року та двох гуманітарних дисциплін. В результаті вдалося вирізнити не лише найактуальніші піднапрями у межах кожної галузі, що дає змогу детальніше описати тематичний спектр чорнобильських досліджень, але й простежити їх в часі, спостерігаючи, як одна тематика змінює іншу.

Дослідження проведено у межах проектів: “Статистична фізика у різноманітних реалізаціях” (Сьома рамкова угода, FP7-PEOPLE, IRSES project N295302) та “Структура та еволюція складних систем із застосуванням у фізиці та природничих науках” (Сьома рамкова угода, FP7-PEOPLE, IRSES project N612669). Особлива подяка колегам по проекту, у який увійшла ця задача: Юрію Головачу, Ральфу Кенні та Бертрану Бершу, а також Нілу ван Еку за плідні дискусії та роз'яснення певних моментів роботи програми VOSviewer.

1. Tseng, Y. H., Lin, Y. I., Lee, Y. Y., Hung, W. C., \ Lee, C. H. (2009). A comparison of methods for detecting hot topics // *Scientometrics*, 81(1), 73–90. 2. Akritidis, L., Katsaros, D., Bozanis, P. (2012). Identifying attractive research fields for new scientists. *Scientometrics*, 91(3), 869–894. 3. Griffith, B. C., Small, H., Stonehill, J. A., Dey, S. (1974). The structure of scientific literatures II: Toward a macro- and microstructure for science. *Science Studies*, 4(4), 339-365. 4. White, H. D., Griffith, B. C. (1981). Author cocitation: A literature measure of intellectual structure. *Journal of the American Society for Information Science*, 32(3), 163-171. 5. Rip, A., Courtial, J. (1984). Co-word maps of biotechnology: An example of cognitive scientometrics. *Scientometrics*, 6(6), 381–400. 6. Tseng, Y. H. (1998, August). Multilingual keyword extraction for term suggestion. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval* (pp. 377–378). ACM. 7. Jones, L. P., Gassie Jr, E. W., Radhakrishnan, S. (1990). INDEX: The statistical basis for an automatic conceptual phrase-indexing system. *Journal of the American Society for Information Science* (1986-1998), 41(2), 87. 8. Kageura, Kyo, and Bin Umino. Methods of automatic term recognition: A review // *Terminology* 3.2 (1996): 259-289. 9. Schneider, J. W. (2005, June). Verification of bibliometric methods' applicability for thesaurus construction. In *ACM SIGIR Forum* (Vol. 39, No. 1, pp. 63-64). ACM. 10. Zipf, G. K. (1949). Human behavior and the principle of least effort. 11. van Eck, N. J. (2011). Methodological advances in bibliometric mapping of science (No. EPS-2011-247-LIS). Erasmus Research Institute of Management (ERIM). 12. Resnick, A. (1961). Relative effectiveness of document titles and abstracts for determining relevance of documents. *Science*, 134(3484), 1004-1006. 13. Zuccala, A., Van Eck, N. J. (2011). Poverty

research in a development policy context. *Development Policy Review*, 29(3), 311-330. 14. Мриглод О. І., Головач Ю. В. (2012). Реакція наукової спільноти на Чорнобильську аварію: аналіз розвитку тематики публікацій // *Вісник НАН України*. 15. Mryglod O., Holovatch Yu., Kenna R., Berche B. *Quantifying the evolution of a scientific topic: reaction of the academic community to the Chornobyl disaster* // *Scientometrics* (подано до друку). 16. Van Eck, N. J., Waltman, L. (2010). *Software survey: VOSviewer, a computer program for bibliometric mapping* // *Scientometrics*, 84(2), 523–538. 17. Van Eck N. J., Waltman L. (2011). *Text mining and visualization using VOSviewer*. arXiv preprint arXiv:1109.2058. 18. *TreeTagger* (2015): a language independent part-of-speech tagger <http://www.cis.uni-muenchen.de/~schmid/tools/TreeTagger/>. Перевірено доступність 15 вересня 2015 р. 19. van Eck, N., Waltman, L., Noyons, E., Buter, R. (2010). *Automatic term identification for bibliometric mapping* // *Scientometrics*, 82(3), 581-596. 20. Simon, H. A. (1955). *On a class of skew distribution functions*. *Biometrika*, 425–440.