

ІНФОРМАЦІЙНІ СИСТЕМИ, МЕРЕЖІ ТА ТЕХНОЛОГІЇ

УДК 004.75

Т. М. Басюк, А. С. Василюк

Національний університет “Львівська політехніка”,
кафедра інформаційних систем та мереж

ФАКТОРИ РАНЖУВАННЯ ІНТЕРНЕТ-РЕСУРСІВ ПОШУКОВОЮ СИСТЕМОЮ GOOGLE

© Басюк Т. М., Василюк А. С., 2016

Описано основні фактори, що впливають на особливості проєктування високорейтингових ресурсів згідно з алгоритмами пошукової системи Google. Проаналізовано відомі технології та методики оцінювання інтернет-ресурсів, що дало змогу виявити фактори, що мають найбільше значення в процесі побудови ресурсу. Наведено методики оцінювання історії ресурсу (archive.org), визначення санкцій за окремим веб-сайтом (Google Panda, Google Penguin), швидкості завантаження з віддалених вузлів й часу доступу (Sitespeed), структури зовнішніх посилань (LinkPad) та підбору ключових слів (Google Adwords Keyword Planner Tool).

Ключові слова: інтернет-ресурс, рейтинг, пошукова система, фактор ранжування.

The paper describes the main factors influencing the design features of high rating resources according to the algorithms search engine Google. The analysis of known technologies and methods of evaluating online resources has been conducted, that made it possible to identify the factors that are the most important in the process of resource building. The methods for the resource history evaluation (archive.org), for determining sanctions imposed on a separate website (Google Panda, Google Penguin), methods for increasing download speed from remote locations and reducing Web access time (Sitespeed), the structure of external links (LinkPad) and keyword selection (Google Adwords Keyword Planner Tool) have been provided.

Key words: Internet resource, ranking, search engine, ranking factor.

Вступ. Загальна постановка проблеми

Сьогодні нікого не здивує наявність інтернет-ресурсу, що виконує різні ролі: від звичайної візитівки (невеликий сайт, що складається з однієї або кількох сторінок та містить основну інформацію про організацію, приватну особу, товари або послуги, контактні дані) до потужного інформаційного порталу (сайт, що містить вичерпний обсяг інформації з будь-якої предметної області або декількох областей, коли основний акцент зроблено на інформаційному наповненні: грамотність подання, зручність читання, подання інформації тощо) чи інтернет-магазину глобального масштабу (сайт, призначений для вибору товарів з каталогу і продажу їх з використанням мережі Інтернет, що зазвичай створюється із застосуванням CMS (Content Management System) та оснащений необхідним функціоналом). При цьому мати свій ресурс у мережі Інтернет не лише престижно, але й прибутково. Аналізуючи відомі фінансові огляди [1], можна зробити висновок, що рейтинг інформаційного ресурсу компанії є співвімірним з показником успішності її діяльності.

Успішність будь-якого інтернет-ресурсу напряму залежить від його популярності серед пошукових систем. При цьому важливо на початку популяризації ресурсу зробити акцент на популярних пошукових системах у цьому регіоні, оскільки саме їх зазвичай обираємо потенційний користувач. Станом на 2015 рік найпопулярнішою на просторах глобальної мережі є пошукова система Google. Вона є лідером у всіх країнах світу, за винятком Росії (Yandex), Китаю (Baidu), Чехії (Seznam) та Північної Кореї (Naver) [2]. Що стосується України, то використання пошуку від компанії Google визначається на рівні 70 %, а його найближчого конкурента Яндекс – 24 %, тобто приблизно втричі менше. Тому для просування інтернет-ресурсів в Україні основний акцент необхідно зробити

саме на пошуковій системі Google, оскільки вона зможе забезпечити найбільшу кількість потенційних користувачів у цьому сегменті ринку.

Зв'язок висвітленої проблеми із важливими науковими та практичними завданнями

Визначення основних факторів ранжування, що використовуються пошуковою системою Google, є надзвичайно важливою та водночас складною задачею. Оскільки, з одного боку, надає необхідний механізм із проектування якісного ресурсу, що вже буде оптимізований й зможе при відповідному розміщенні входити в ТОП пошукової системи, а з іншого – оскільки цей продукт є комерційною розробкою то визначення його основних алгоритмів оцінювання полягатиме в аналізі відповідних ресурсів та позицій, які вони займають, що є надзвичайною трудомісткою та протяжною в часі роботою. З огляду на те, визначена задача являє собою складний та безперервний процес, під час здійснення якого необхідно застосовувати розширеній математичний та алгоритмічний апарат. Зокрема можна застосувати як алгоритми пошуку інформації, так і методи оцінювання рейтингу ресурсу та його визначення згідно з показником PR (Page Rank). Розв'язання поставленої задачі дасть змогу визначити основні фактори ранжування інтернет-ресурсу, сприятиме побудові оптимізованих високорейтингових ресурсів та надасть засоби із адаптації відомих методів та алгоритмів до задач побудови інтелектуальних інформаційних систем у цій предметній галузі.

Аналіз останніх досліджень та публікацій

Як було показано [3], сьогодні існує множина методик із визначення рейтингу інтернет-ресурсу, починаючи від локальних алгоритмів, які використовуються пошуковими системами (PR, TiЦ), до глобальних, що являють собою сукупність множини показників. Рейтинг інтернет-ресурсу, а відтак і відповідна його позиція визначається множиною факторів, які фактично є малодослідженими внаслідок закритості пошукової системи та авторських прав. Наявні сьогодні дослідження [4, 5] у загальному вигляді описують особливості роботи пошукових систем та можливі фактори, які впливають на популяризацію ресурсу, оминаючи при цьому необхідні механізми ранжування, що використовуються найпопулярнішою пошуковою системою. З огляду на те роботу із створення ресурсу в більшості випадків ділять на дві частини. Перша – власне розроблення інтернет-ресурсу відповідно до критеріїв користувача та розміщення в глобальній мережі, друга – його популяризація (SEO) під пошукові системи. Зазначений алгоритм роботи вимагає як значних обсягів часу (з огляду на тривалість процесу), так і фінансових витрат (фактично виконується одна й та сама робота, проте за подвійну плату).

Особливістю проведених досліджень є те, що визначення факторів, які впливають на позицію інтернет-ресурсу в пошуковій видачі Google, робить актуальною задачу проектування системи оцінювання та надання рекомендацій із створення веб-сайтів, що будуть оптимізовані під відповідні запити користувачів у мережі Інтернет.

Основні завдання дослідження та їх значення

Метою дослідження є визначення основних факторів ранжування інтернет-ресурсів у пошуковій системі Google. Проведене дослідження надасть засоби із початкової побудови веб-сайту згідно з оптимізаційними чинниками, що сприятиме його входженню у топові позиції рейтингу пошукової системи Google. Для досягнення поставленої мети необхідно вирішити такі основні завдання: проаналізувати відомі технології та методики оцінювання інтернет-ресурсів інформаційно-пошуковою системою Google; виявити фактори, що найбільше впливають на процес побудови високорейтингового ресурсу; проаналізувати основні методики оцінювання історії ресурсу, визначення санкцій за окремим веб-сайтом, швидкості завантаження з віддалених вузлів й часу доступу, структури зовнішніх посилань та підбору ключових слів.

Спроектований з використанням цих факторів інтернет-ресурс зможе претендувати на топові позиції в рейтингу як пошукової системи Google, так й інших, які працюють за схожими принципами, що, своєю чергою, призведе до його популяризації в глобальній мережі. Результати виконаної роботи у своїй сукупності розв'язують актуальну задачу визначення факторів ранжування інтернет-ресурсу, що надасть засоби із його популяризації в глобальній мережі.

Основні результати досліджень

Історія пошукового гіганту розпочалась в березні 1998 року у вигляді наукового проекту Л. Пейджа та С. Бріна – тоді студентів Стендфордського університету. Вони працювали над ідеєю створення перспективної технології побудови єдиної інтегрованої цифрової бібліотеки. Разом з тим Пейдж вивчав математичні властивості глобальної мережі, представляючи структуру посилального оточення веб-сайту у вигляді графу [4]. При такому представленні вдалося вивчити характер посилань з погляду цінності інформації на сторінці (за аналогією із цитатами в наукових виданнях). Перетворення отриманих даних сприяло виникненню принципово нового критерію оцінювання, що отримав назву PageRank (PR) [5]. Особливістю його є те, що при формуванні списку пошукової видачі розміщення сторінки залежить від кількості посилань з інших сайтів на неї та від їх популярності.

При цьому PageRank – це нормалізоване відношення кількості посилань, що надходять на сторінку, до кількості посилань, що виходять з неї. Розрахункова формула має вигляд:

$$PR = (1 - d) + d(PR(T_1)/C(T_1)) + \dots PR(T_n)/C(T_n) \quad (1)$$

де d – емпірично підібраний коефіцієнт ($d=0.85$); $T_1 \dots T_n$ – інтернет-ресурси, які містять посилання на документ; $C(T_1) \dots C(T_n)$ – загальна кількість посилань, які надходять із сторінок $T_1 \dots T_n$.

З формули (1) випливає, що PageRank будь-якого ресурсу залежить від PageRank ресурсів, з яких можливий перехід на нього. Отже, цей показник буде завжди високим для веб-сторінок, що є популярними в глобальній мережі. PageRank має значну ймовірнісну складову, яка визначає потрапляння довільного користувача на визначену сторінку. Сума зазначеного показника для всіх сторінок дорівнює одиниці:

$$\sum_j PR_j = 1, j = 1 \dots N, \quad (2)$$

де N – кількість проіндексованих сторінок.

Можна зауважити (2), що у випадку штучного збільшення популярності інтернет-ресурсу достатньо побудувати деяку кількість взаємопов'язаних веб-сторінок і так збільшити PageRank. Зокрема створити велике кільце веб-сайтів, в яких інтернет-ресурс містить лише “дружні посилання”. Тоді кожен з них внаслідок ітеративності алгоритмів визначення PageRank матиме високий коефіцієнт “важливості”. Незважаючи на це припущення, в моделі функціонування пошукової системи Google така ситуація не спрацьовує внаслідок застосування запатентованих механізмів, що не дозволяють нараховувати занадто високий PageRank “підозрілим” ресурсам.

Іншою важливою особливістю пошукової системи є збереження опису посилань на проіндексовані веб-сторінки. Зокрема, якщо автор сторінки забув вказати назву між тегами `<title>` `</title>` (будь-яка система при пошуковій видачі ставить найвищий пріоритет назві), система Google буде орієнтуватися за текстами посилань на неї, ґрунтуючись на припущеннях, що у випадку розміщення посилань на веб-сторінку автор останньої її вивчив і постарається якнайповніше відобразити зміст у тексті посилання. Зазначена особливість надає гнучкі механізми пошуку в накопичений базі даних системи.

Крім розрахунків PageRank і запам'ятовування тексту посилань, пошукова система зберігає шрифтовий розмір та зсув кожного слова відносно початку документа. З огляду на те, в результаті пошукового запиту можна вивести документ, в якому це слово виділено або розміром шрифта, або знаходиться поблизу початку документа. Крім того, оскільки система володіє інформацією щодо місця кожного слова, на практиці стає можливий так званий Proximity search (пошук за найближчим розташуванням слів). Наприклад, за запитом `<СЛОВО1 СЛОВО2>` пошукова система знайде множину документів, проте у верхніх позиціях розташуються ті, в яких `<СЛОВО1>` знаходиться максимально близько від `<СЛОВА2>`.

На основі проведених досліджень щодо позиції інтернет-ресурсу в результатах пошукової видачі системи Google вдалося виявити такі основні фактори ранжування (рис. 1). Найважливішим фактором, що впливає на просування інтернет-ресурсу, є доменне ім'я (назва, на яку будуть розміщувати посилання інші ресурси та яка буде відображена у пошуковій видачі). На нього найбільше впливають: вік, історія, ключові слова, адреса, субдомен та історія санкцій.



Рис. 1. Фактори ранжування пошукової системи Google

Вік домену відіграє значну роль, оскільки в процесі просування нового сайту (вік до шести місяців) необхідно здійснити множину перетворень, починаючи з внутрішньої структури, завершуючи зовнішнім посилальним оточенням. Що стосується внутрішньої структури, то насамперед необхідно враховувати такі фактори [6]:

- *метатеги* (Title – основний метатег, що враховується при ранжуванні ресурсу; Description – застосовується для підвищення індексу CTR за результатами пошуку; Keywords – визначає ключові слова ресурсу);
- *теги* (ALT – призначений для опису зображень і визначає альтернативний текст, його заповнення та сприяє реалізації пошуку за графічними зображеннями; H1 – визначає назву сторінки та застосовується не більше одного разу; Текст – власне текстове наповнення);
- *перелінковка інтернет-ресурсу* – надзвичайно важливий фактор, який необхідно враховувати при внутрішній оптимізації, особливо у випадку наявності множини веб-сторінок. Найпоширенішими способами є: перелінковка блоками (використовуються спеціальні скрипти, що здійснюють розподіл ваги ключових слів за інтернет-ресурсом), перелінковка в статтях (коли всі сторінки між собою перелінковані та передають ваги [3]);
- *соціальні фактори* – в сучасному інформаційному суспільстві все більше впливають на просування інтернет-ресурсу; з огляду на це наявність кнопок соціальних мереж є обов’язковою на сторінки веб-сайту;
- *юзабіліті* – визначає дизайн та зручність користування ресурсом й значно впливає на: перегляди сторінок, замовлення, продажі. Що зручніший інтернет-ресурс та кращі його структура й сприйняття, то більшу кількість трафіку він залучатиме;
- *код* – необхідно оптимізовувати насамперед, оскільки це безпосередньо впливає на швидкість роботи ресурсу, а відтак на його ранжування. При повільному завантаженні ресурсу пошукова система понижав його позиції при пошуку за ключовими словами;
- *інструменти пошукових систем* – необхідно обов’язково додавати ресурс у Google Webmaster Tools з метою його попереднього аналізу пошуковою системою. Крім того, за допомогою цього інструментарію можна налаштовувати файл robots.txt, що є важливим для внутрішньої оптимізації.

Історія домену може або допомогти, або нашкодити просуванню, наприклад, якщо тематика ресурсу не змінювалась, то це позитивно впливає на ранжування і навпаки. З огляду на це множина SEO-оптимізаторів шукає певні тематичні ресурси з метою створення відповідних веб-сайтів. Історію будь-якого інтернет-ресурсу можна знайти у глобальному архіві ARCHIVE.ORG (<http://archive.org/web/web.php>), в якому індексуються веб-сайти глобальної мережі у різні проміжки часу. Зокрема, історію ресурсу Національного університету “Львівська

політехніка” збережено від 2003 року. Отже, можна вибрати рік, наприклад, (25.06.2003 р.) та переглянути вигляд ресурсу в цей час (рис. 2).

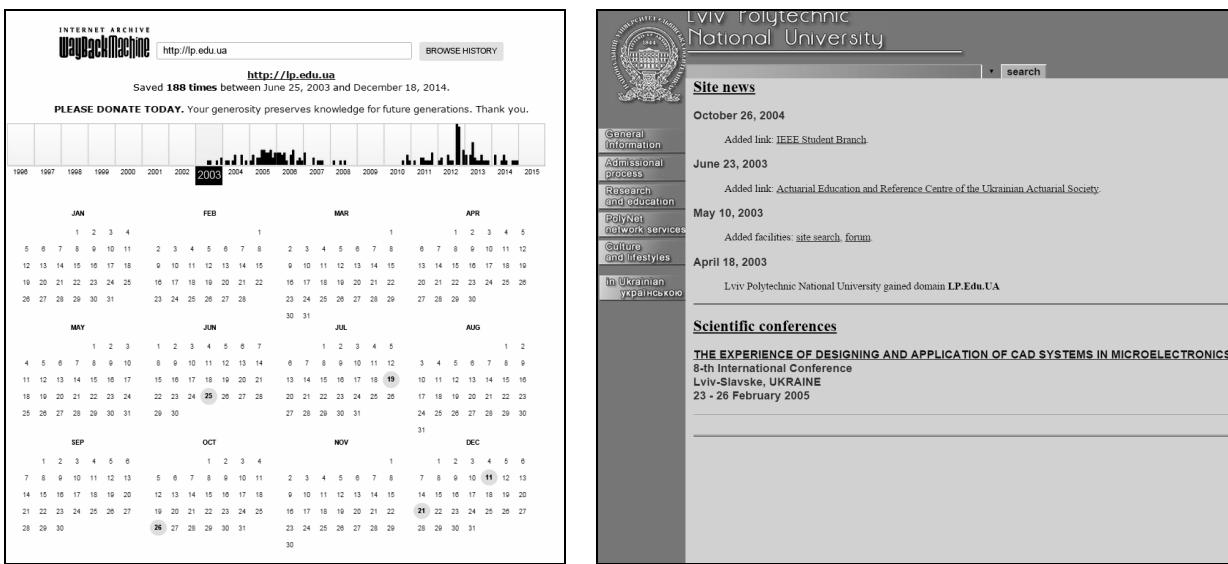


Рис. 2. Історія ресурсу Національного університету “Львівська політехніка”

Ключові слова в домені дають відчутний результат у процесі просування, проте з введенням алгоритму EMD (Exact Match Domains) домени, які містять у назві множину ключових слів або пошукових запитів, свідомо блокуються.

Субдомен – практика використання субдоменів виникла в результаті використання об’ємних інтернет-ресурсів, оскільки розділи різного типу простіше окремо просувати та модифікувати.

Історія санкцій домену суттєво впливає на процес просування ресурсу, оскільки пошукова система оцінює також і минулі санкції. Google застосовує множину алгоритмів, серед поширених – Google Panda та Google Penguin. Алгоритм Google Panda, запроваджений 23 лютого 2011 року, має на меті підвищення якості результатів пошуку оцінюванням матеріалів ресурсу. Основними напрямками роботи алгоритму є:

- *недостатній контент* – передбачається, що інформаційне наповнення ресурсу не є новим або цінним для користувача чи не розкриває належно інформаційної теми;
- *дубльований контент* – визначає використання копій контенту з інших джерел;
- *некласичний контент* – пошукова система Google залучає інтернет-ресурси, які постійно оновлюються, з огляду на те оптимізатори з метою оновлення сторінок можуть публікувати некласичний контент, що, своєю чергою, принесе більше шкоди, ніж користі.

Google Penguin введено до пошукової системи 24 квітня 2012 року для боротьби з інтернет-ресурсами, які застосовують сумнівні зворотні посилання. Авторитет та вага інтернет-ресурсу значною мірою залежать від того, які веб-сайти на нього посилаються. При цьому одне посилання з авторитетного джерела може мати таку саму вагу, як десятки посилань з маловідомих ресурсів. Робота алгоритму Google Penguin полягає у виявленні маніпуляцій з посиланнями, які створені оптимізаторами вручну з метою вплинути на позиції інтернет-ресурсу. До таких засобів належать: обмін посиланнями; посилання з нерелевантних ресурсів, участь у договірних схемах, посилання із сайтів-сателітів тощо.

Іншим фактором, який впливає на ранжування ресурсу, є *сервер (хостинг)*. Основними критеріями, які можуть вплинути, є: довгий відклик інтернет-ресурсу, недоступність інтернет-ресурсу (особливо коли Google Bot намагається зайти, щоб проіндексувати ресурс), Uptime ресурсу (час неперервної роботи обчислювальної системи чи її частини, вимірюється з моменту завантаження до моменту зупинки роботи чи завершення моніторингу). Цей фактор можна проаналізувати за допомогою інтернет-ресурсу <http://sitespeed.me/>, в якому відображається інформація щодо швидкості завантаження веб-сторінок із різних країн світу (рис. 3).

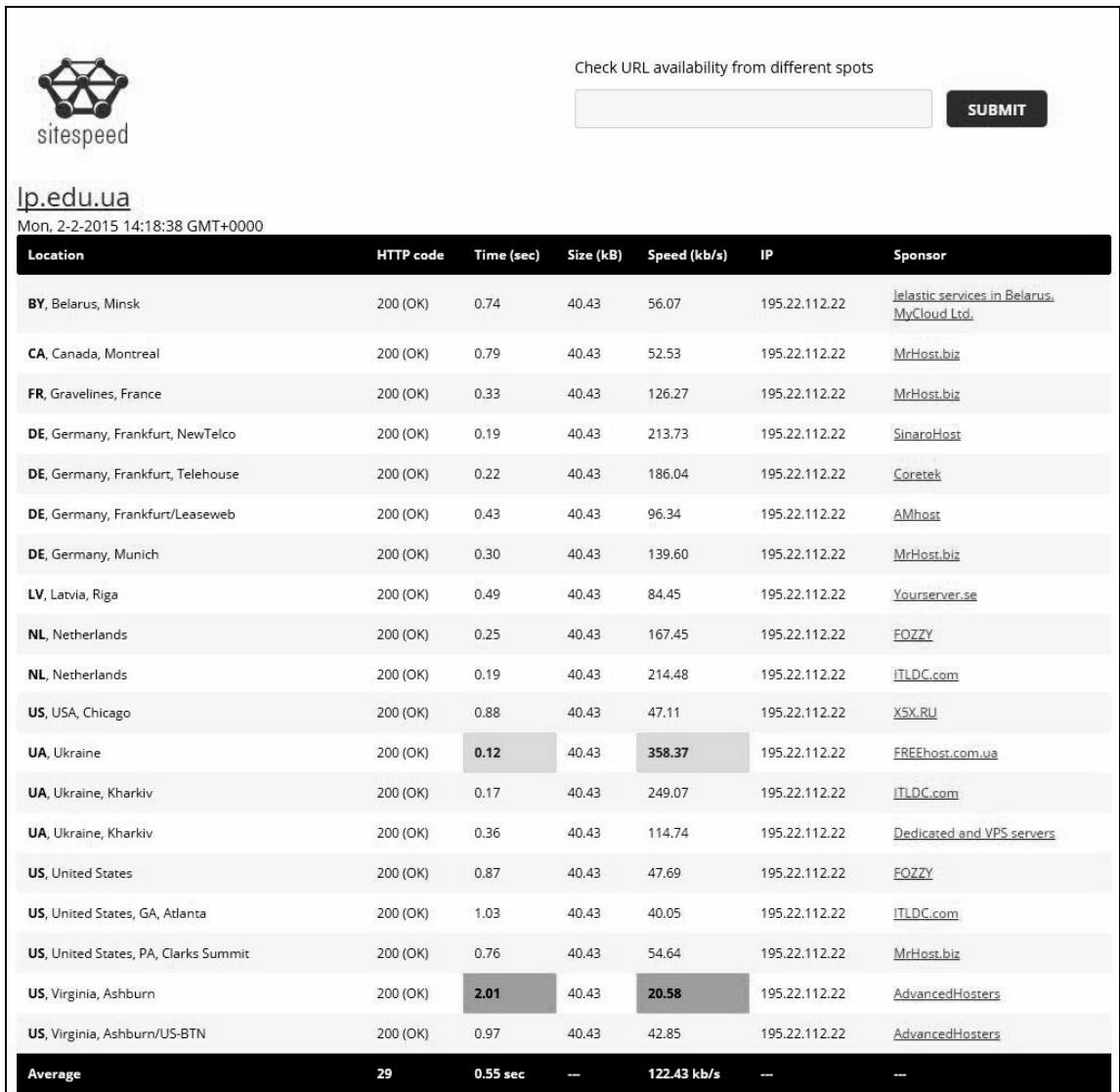


Рис. 3. Швидкість завантаження ресурсу lp.edu.ua

Результатуюча інформація відображається у вигляді середнього значення (average) за часом доступу, розміром отриманої інформації та швидкістю доступу.

Контент – основний елемент інтернет-ресурсу, який підлягає обов’язковій оптимізації з метою підвищення популяризації ресурсу. Все наповнення повинно відповідати певній структурі (рівень вкладеності, відсутність дубльованого контенту, наявність мікроданих) та бути інформативним для користувача.

Існує два види вихідних посилань, що впливають на просування інтернет-ресурсу: *внутрішні* (ведуть на сторінки веб-сайту) та *зовнішні* (направлені на інші ресурси). При цьому з використанням внутрішніх посилань можна рівномірно розподілити їх вагу між сторінками інтернет-ресурсу, що дає змогу просувати веб-сайт не лише за загальними запитами (головної сторінки), але і за більш спеціалізованими – внутрішніх розділів.

Зовнішні посилання на інші ресурси оцінюють за допомогою технології TrustRank, що є фільтром пошукової системи. Завданням TrustRank є відділення “корисних” ресурсів від “зайвих”, створених виключно для просування інтернет-ресурсу. Зокрема, якщо посилань із ресурсу багато і вони ведуть на веб-сайти з низькою відвідуваністю і авторитетністю, TrustRank запускає алгоритм зниження ваги ресурсу у пошуковій видачі.

З метою ранжування ресурсу необхідно дотримуватись таких критеріїв [7]: кількість посилань (оптимальна кількість вихідних посилань залежить від тематики інтернет-ресурсу, причому на одній сторінці їх не може бути більше п’яти); тематичність (інтернет-ресурс, зі сторінок

якого йдуть посилання на ресурси різної тематики, може бути сприйнятій пошуковою системою як рекламний простір з продажу посилань, що призводить до зниження рейтингу); авторитетність (вихідні посилання на ресурси з високими Page Rank позитивно оцінюються пошуковою системою, особливо, якщо їх авторитетність вища, ніж ресурсу, який просувається). При цьому на просторах глобальної мережі існує множина ресурсів, які надають вичерпну інформацію щодо кількості внутрішніх та зовнішніх посилань. Одним з найпоширеніших є ресурс: <https://www.linkpad.ru/> (рис. 4).

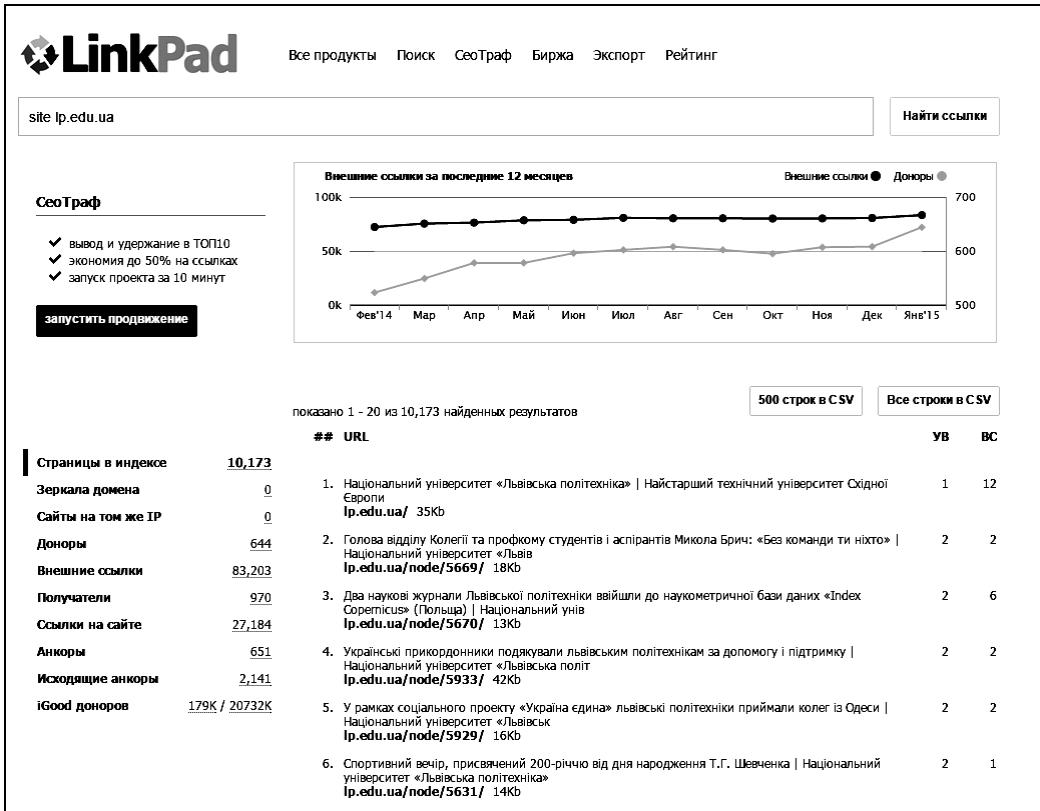


Рис. 4. Інформація про ресурс lp.edu.ua

Зворотні посилання – це посилання, які вказують на інтернет-ресурс користувача і є одним з найважливіших факторів пошукової оптимізації, а їх аналіз – пріоритетний чинник у процесі ранжування ресурсу. Отримання якісних зворотних посилань збільшує цільовий трафік, а відтак сприяє входженню інтернет-ресурсу до топових позицій пошукової видачі. Існує два типи зворотних посилань: односторонні (містять посилання на ресурс користувача) та двосторонні (містять взаємні посилання між ресурсами). Перший тип значною мірою впливає на рейтинг ресурсу та переважно визначається такими чинниками посилального оточення: якість інтернет-ресурсу та сторінки посилання, вік ресурсу та сторінки, ролевантність контенту сторінки, розташування посилання та анкор посилання. Другий тип майже не враховується пошуковою системою Google.

Пошукові фактори – це додаткові користувачькі особливості, що впливають на ранжування ресурсу. Серед найпоширеніших чинників є: показник відмов та різне сприйняття інформації. Під показником відмов розуміють відношення кількості відмов до кількості відвідувачів ресурсу. Формально система Google Analytics *відмовою* вважає ситуацію, за якої в процесі відвідування інтернет-ресурсу, крім перегляду сторінки, не було зафіксовано інших дій. Цей показник слугує своєрідним критерієм якості аудиторії або сторінки інтернет-ресурсу: що вищий показник, то гірша якість. Адже якщо більша частина відвідувачів заходять та залишають ресурс, нічого не зробивши, вони або не зацікавлені у його тематиці (нецільова аудиторія), або ця веб-сторінка не сподобалася (незрозуміла або не викликає довіри). З огляду на це потрібно: змінити оформлення сторінок входу і цільових сторінок; змінити сторінки так, щоб вони якомога точніше відповідали пошуковим запитам користувачів; змінити ключові слова, щоб вони точніше відповідали вмісту сторінок [8].

З метою ефективного підбору ключових слів підзапитів, що надходять з пошукової системи Google, можна використовувати утиліту Google Adwords Keyword Planner Tool (рис. 5). Особливістю цього сервісу є, окрім підбору ключових слів ресурсу, підбір контекстної реклами та оцінювання трафіку (кількості користувачів, котрі можуть потрапити на веб-сайт).

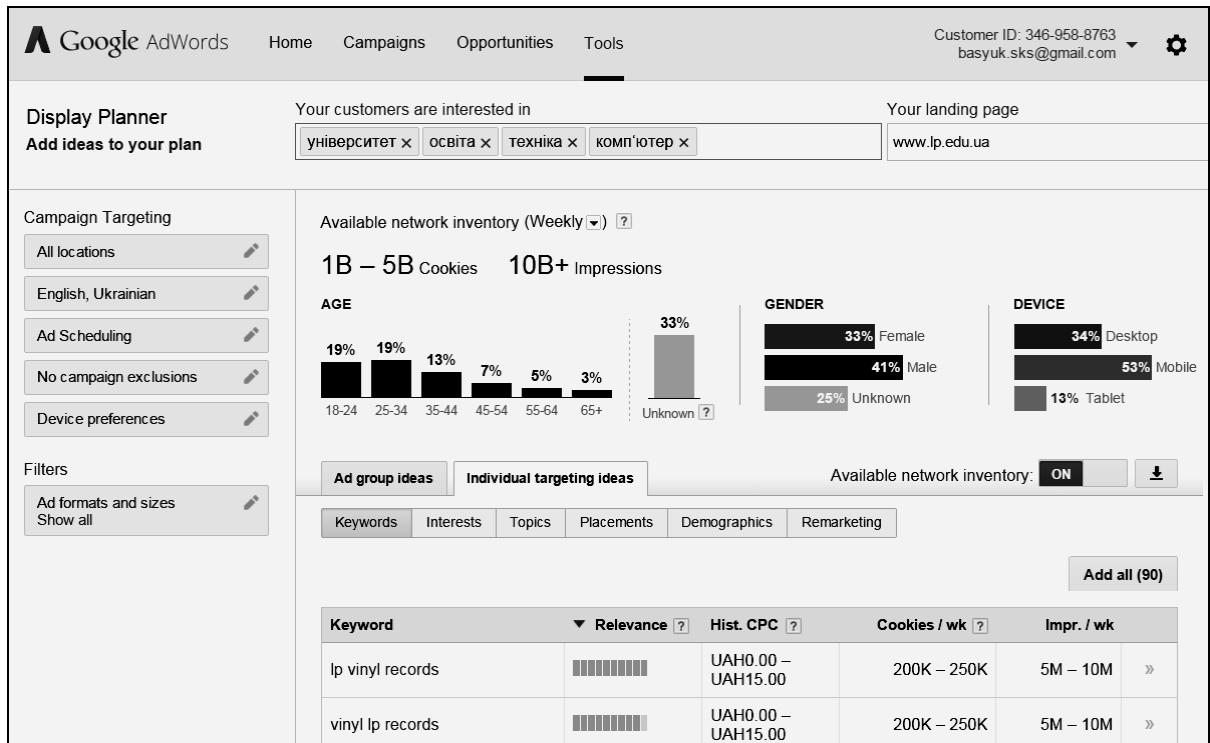


Рис. 5. Вікно сервісу Google Adwords

Висновки та перспективи подальших наукових розвідок

У результаті проведеного дослідження виявилось, що пошукова система Google містить значну кількість (понад п'ятсот) факторів ранжування. При цьому для успішної побудови високорейтингових ресурсів необхідно впливати на ті, які є відомими, з метою отримання вагомого результату у вигляді постійного росту відвідуваності веб-сайту. Тому необхідною умовою є створення якісного, цікавого та корисного контенту, який надасть засоби із вирішення завдань користувачів – тоді перші місця в пошукових видачах будуть забезпечені.

Подальші дослідження будуть направлені на створення критеріїв оцінювання інтернет-ресурсів, що дасть змогу визначити рівень підготовки веб-сайту та можливості його ранжування.

1. Стрикчиола Д. SEO. Искусство раскрутки сайтов / Джесси Стрикчиола, Рэнд Фишкин, Стефан Спенсер. – СПб.: ВНВ, 2011. – 520 с.
2. Лоу Д. Google. Прошлое, настоящее, будущее / Джанет Лоу. – М.: Манн, Иванов и Фербер, 2009. – 320 с.
3. Басюк Т. М. Принципы побудови системи аналізу та просування інтернет ресурсів / Т. М. Басюк // Комп'ютерні науки та інформаційні технології. – Львів: Нац. ун-т “Львівська політехніка”. – 2012. – № 784. – С. 43–48.
4. Граппоне Д. Поисковая оптимизация сайтов. Исчерпывающее руководство / Дженнifer Граппоне, Градива Казн. – М.: Эксмо, 2012. – 528 с.
5. Ледфорд Д. SEO / Джерри Ледфорд. – М.: Экран, 2009. – 304 с.
6. Басюк Т. М. Проектування системи автоматизації процесу seo-оптимізації / Т. М. Басюк // Комп'ютерні науки та інформаційні технології. – Львів: Нац. ун-т “Львівська політехніка”. – 2014. – № 800. – С. 92–87.
7. Basyuk T. M. Influence of external factors on the position rankings website / T. M. Basyuk // Тези доповідей третьої міжнародної науково-практичної конференції “Інформаційні управлюючі системи та технології”. – Одеса: BMB, 2014. – С. 241–242.
8. Геддс Б. Google AdWords. Исчерпывающее руководство / Брэд Геддс. – М.: Манн, Иванов и Фербер, 2014. – 950 с.