


Promoting Document Relevance Using Query Term Proximity for Exploratory Search

Vikram Singh, National Institute of Technology, Kurukshetra, India*

 <https://orcid.org/0000-0001-6315-0872>

ABSTRACT

In the information retrieval system, relevance manifestation is pivotal and regularly based on document-term statistics, i.e., term frequency (tf), inverse document frequency (idf), etc. Query term proximity (QTP) within matched documents is mostly under-explored. In this article, a novel information retrieval framework is proposed to promote the documents among all relevant retrieved ones. The relevance estimation is a weighted combination of document statistics and query term statistics, and term-term proximity is simply aggregates of diverse user preferences aspects in query formation, thus adapted into the framework with conventional relevance measures. Intuitively, QTP is exploited to promote the documents for balanced exploitation-exploration, and eventually navigate a search towards goals. The evaluation asserts the usability of QTP measures to balance several seeking tradeoffs, e.g., relevance, novelty, result diversification (coverage, topicality), and overall retrieval. The assessment of user search trails indicates significant growth in a learning outcome (due to novelty).

KEYWORDS

Exploratory Search, Information Retrieval, Query Term Proximity, Relevance, Retrieval Strategy

INTRODUCTION

Information-seeking is a fundamental endeavor of human being and several information search systems has been designed to assist a user to pose queries and retrieves informative data to accomplish search goals. The traditional systems strongly trust user's capability of phrasing precise request and perform better if requests are short and navigational. A potential obstacle to such systems is an astonishing rate of information overload that makes difficult to a user for identifying useful information. Therefore nowadays, search focus is shifting from finding to understanding information (White & Roth, 2009), especially in *discovery-oriented* search. When a user wants information for learning purpose, decision making or other cognitive activity, the conventional search methodologies are not capable to assist, though data exploration is helpful. A data exploration synthesis *focused* search

DOI: 10.4018/IJIRR.325072

*Corresponding Author

This article published as an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>) which permits unrestricted use, distribution, and production in any medium, provided the author of the original work and original publication source are properly credited.

and *exploratory* browsing, to discover the interesting data objects. Though, exploration become a recall-oriented navigation over complex and huge datasets using short typed ill-phrased data request (Idreos, Papaemmanouil, & Chaudhuri, 2015; White, 2016; Marchionini, 2006), and thus requires strong support for adaptive relevance measures in retrieval framework (Nandi, & Jagadish, 2011).

In the data deluge, retrieval of relevant data requires either formal awareness of complex schema and content for the formulation of a data retrieval request or assistance from information system (Kersten, Idreos, Manegold, & Liarou, 2011; Huston, Culpepper, & Croft, 2014). For both situations, the system employs *implicit* measures to outline matched objects and *explicit* measures to eventually steer search towards a *region-of-interest*. Most existing retrieval models score a document predominantly on documents-terms statistics, i.e. *document lengths*, *query-term frequencies*, *inverse document frequencies*, etc (Van, 1977; Daoud & Huang, 2013). Intuitively, the *query terms proximities* (QTPs) within pre-fetched result set/documents could be exploited for re-position/re-raking of the documents/results in which the matched query terms are close to each other. For example, an information search considering the query '*exploratory search*' on two documents, both matching the two query terms once:

$Doc_1: \{...exploratory search.....\}.$

$Doc_2: \{....exploratory....search....\}.$

Intuitively, *document₁* should be ranked higher, as occurrences of both query terms are closest to each other. In compare to the *document₂*, where both query terms are far apart and their combination does not necessarily imply the meaning of '*exploratory search*'.

The *term-term* affinity within matched document has role to play during the retrieval and eventually to position the document in appropriate relevance (Salton & Buckley, 1988; Borlund, 2003; Verma, 2016). For an information search, a user specify data request in more than one terms with an anticipated inherent closeness. The closeness in query terms characterizes structural constraints of a user query and the importance between two matched documents in an information-seeking. The *query term proximity* is one measure, however, has been principally under-explored in traditional retrieval framework and models; mainly due to intrinsic design concerns (*how we can model proximity*) and its overall usability (*what it serve*) into a retrieval model.

This paper systematically explores the query term proximity heuristic, to guide the user's information-seeking by deliberating both *document-terms* (DTs) and *query-terms* (QTs) relevance means. The focus is on three *research questions* (RQs):

RQ1: *What constitutes relevance in exploitation and exploration? What relevance type is most significant?*

RQ2: *How can query term proximity (QTP) be adapted with document-terms relevance to optimize information exploitation and eventually exploration efforts?*

RQ3: *Finally, how to design an information retrieval framework that, account user's search task while rewarding or penalizing both relevance measures.*

For a given *document corpus* and *query terms* (in user's query), relevance manifestation is done across three factors: coincidence of QTs with DTs (*intra-document* and *inter-document* relevance) and *Span* of QTs (*intra-document*) and distance of QTs (*intra-document*). The significance of measures is derived via a study on user defined *search trails* (STs), and eventually the impact of overall retrieval framework on the exploration efforts.

20 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/article/promoting-document-relevance-using-query-term-proximity-for-exploratory-search/325072

Related Content

Personalized Content-Based Image Retrieval

Iker Gondra (2008). *Personalized Information Retrieval and Access: Concepts, Methods and Practices* (pp. 194-219).

www.irma-international.org/chapter/personalized-content-based-image-retrieval/28074

Service-Driven Computing: Challenges and Trends

Raja Ramanathan (2018). *Information Retrieval and Management: Concepts, Methodologies, Tools, and Applications* (pp. 2287-2311).

www.irma-international.org/chapter/service-driven-computing/198648

Efficacious Hyperlink Based Similarity Measure Using Heterogeneous Propagation of PageRank Scores

Vasantha Thangasamy (2019). *International Journal of Information Retrieval Research* (pp. 36-49).

www.irma-international.org/article/efficacious-hyperlink-based-similarity-measure-using-heterogeneous-propagation-of-pagerank-scores/236655

Tweet Sentiment Analysis with Latent Dirichlet Allocation

Masahiro Ohmura, Koh Kakusho and Takeshi Okadome (2014). *International Journal of Information Retrieval Research* (pp. 66-79).

www.irma-international.org/article/tweet-sentiment-analysis-with-latent-dirichlet-allocation/127002

Emotion Recognition Using Facial Expressions

Arush Jasuja and Sonia Rathee (2021). *International Journal of Information Retrieval Research* (pp. 1-17).

www.irma-international.org/article/emotion-recognition-using-facial-expressions/280523