

Consecutive Time Sensitive Queries Using TF-IDF

M. Madhavi^{#1}, CH. Raja Jacob^{#2}

#1Nova College of Engineering & Technology, Vegavaram, Jangareddy Gudem,
#2M-Tech, CSE, Nova Nova College of Engineering & Technology,
Vegavaram. Jangareddy Gudem

Abstract: Search Engines uses an approach that might help to identify how relevant the results it displays to searchers might actually be, and how likely those results are to show a variety of results when a searcher uses a query term that might cover a range of topics in future. For an important class of queries termed time-sensitive queries over frequently updated archives such as news archives, topic similarity alone is not sufficient for ranking. For such queries, the publication time of the documents is important and should be considered in conjunction with the topic similarity to derive the final document ranking. For incorporating the time dimension, prior systems used an estimation algorithm that considers publication date and time of the documents to locate time periods of interest. However, a document published on the same context at a later date (e.g., a review article, summarizing an event) may also be relevant; We propose to infer the temporal relevance of a document by analyzing its contents, and not by relying solely on its publication date thus increasing the relevancy of the results. So we propose to use Tf-idf, term frequency-inverse document frequency a numerical statistic method, that reflects how important a word is to a document in a collection or corpus. We emulate the performance of the estimation algorithm in combination with tf-idf weights for detecting the important time intervals for a query over a news archive and for incorporating this information in the retrieval process. We show that our techniques are robust and significantly improve result quality for time-sensitive queries compared to state-of-the-art retrieval techniques.

1. INTRODUCTION

TIME is an important dimension of relevance for a large number of searches. Research on searching over such collections has largely focused on retrieving topically similar documents for a query. For a large family of Queries the time dimension are ignoring or not fully exploiting can be detrimental. We should consider not only the document topical relevance but the publication time of the documents as well. There are two motivational points on searching over news archives.

1. Topic-similarity ranking does not model time explicitly, which means that the important dimension of time. But the dimension of time is not considered directly when deciding on the results that are returned for a user query.
2. A topic-similarity ranking of the query results often does not reflect the distribution of relevant documents over time. or many queries, users have a general about the relevant time periods for the queries.

In our schema for an important class of queries over news archives that we call time-sensitive queries. For such queries, the publication time of the documents is important and should be considered in conjunction with the topic similarity to derive the final document ranking. Searching over the large archives highly used method is of timed documents incorporate time in a relatively crude manner: users can submit a keyword query or alternatively sort the results on the publication date of the articles. But, searchers do not always know the appropriate time intervals for their queries, and placing the burden on the users to explicitly handle time during querying is not desirable. Temporal distribution of matching articles for a query, Google's News Archive Search supplements query results with a "timeline". Google also highlights key time periods for each query, so users can explicitly restrict the search to a specific time period. To exploit query result timelines to decide whether to ask users to select appropriate time periods for their queries. We several techniques to estimate the temporal relevance of a day to a query at hand. To estimation the techniques, we use the temporal distribution of matching articles for the query to compute the probability that a day in the archive has a relevant document for the query. Li and

Croft's[2] time-sensitive approach processes a recency query by computing traditional topic-similarity scores for each document, and then "boosts" the scores of the most recent documents, to privilege recent articles over older ones.

Comparing to the traditional models, assume a uniform prior probability of relevance $p(d)$ for each document d in a collection, Li and Croft define the prior $p(d)$ to be a function of document d 's creation date. The the respective to the time prior probability $p(d)$ decreases exponentially. We designed for queries that are after recent documents but the other type of time-sensitive queries are not handled. In our schema we propose a more general framework for answering time-sensitive queries that builds on and substantially expands the earlier work on recency queries.

One alternative is to automatically suggest, based on the query terms, relevant time ranges for the query and allow users to explicitly select appropriate time intervals [3]. As the less input is demanded from the user. But we can automate the previous procedure and prioritize results from periods that we automatically identify as relevant.

Specifically, we design general framework to incorporate time into the retrieval task in a principled manner. These intervals are then used to adjust the document relevance scores by boosting the scores of documents published within the important intervals. Our system provides a web interface for searching the Newsblaster archive, an operational news archive and summarization system, and for experimenting with variations of our approach[4]. We present an extensive evaluation of our system, using both TREC data and real web data analyzed using the Amazon Mechanical Turk[5]. Te resultant show that the quality of the results produced by our techniques for timesensitive queries is significantly higher than that of the (strong) baselines that we consider. We introduce the notion of temporal Relevance which is the probability of a document published at a certain time to be relevant to a given query. We integrate temporal relevance with state-of-the-art retrieval models, including a query likelihood (QL) model, a relevance model (RM), a probabilistic relevance model (PRM), and a query expansion with pseudo relevance feedback model, to naturally process time-sensitive queries.

2. TIME-SENSITIVE QUERIES

The relevant documents may be distributed differently over the time span of a news archive [3]. The relevant results for some queries may exist in certain time periods, large-scale news coverage relevant to the queries takes place and diminishes after a period of time. To illustrate the difference between time-sensitive and time-insensitive queries shown in below figures.

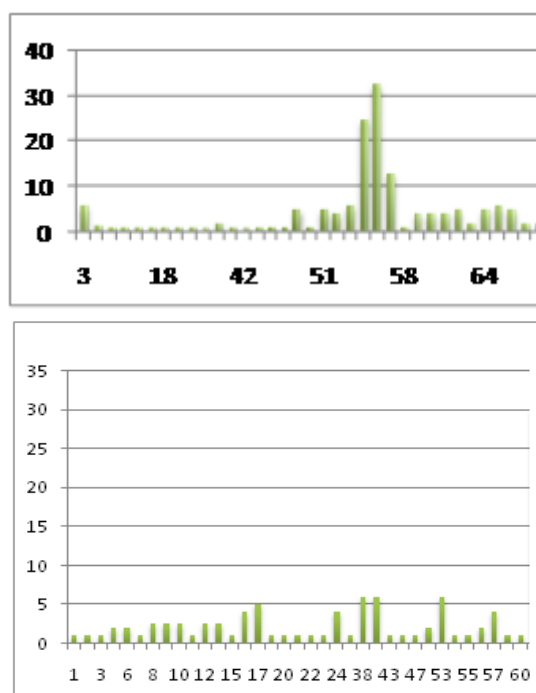


Fig. 1. Relevant-document histograms of a time-sensitive (a) and a timeinsensitive (b) query from a TREC ad hoc query set. (a) Query #311, [Industrial Espionage], a time-sensitive query. (b) Query #304, [Endangered Species (Mammals)], a time- insensitive query.

The figure 1 shows the histogram of both a time-sensitive query (TREC query number 311) and a timeinsensitive query (TREC query number 304). News archives often include many matching documents for time-sensitive queries. Consider the query has 936 matching stories in The New York Times archive, as of March 2009. We said that the traditional topic-similarity ranking alone may not be desirable for time-sensitive queries. Our basic intuition is that the relevance of one document for a given query provides us with useful information about the relevancy of other documents with similar content that were published around the same time. We discuss our first step in accounting for time by introducing techniques to estimate temporal relevance, which is the probability that a time period is relevant to a query at hand.

3. TEMPORAL RELEVANCE

Documents in archival collections are stamped with their publication dates. Queries often are answered and ranked without consideration of these time stamps. To answer the type of time-sensitive queries over a news archive, we would like to use the temporal information implicitly available in the archive. For this, we observe that time-sensitive queries are generally after documents from specific time periods. the majority of documents relevant to the TREC query are located within a specific time period as shown in the figure 2. Hence the distribution of relevant documents over time for a given query. The distribution, hypothesize, can be used to improve the answer quality for a retrieval task. Unfortunately, we usually have no a priori knowledge of the relevant documents for a given query, Hence we cannot accurately compute this distribution.

By Ground Truth

For a given query q , its complete set of relevant documents R_q , which we refer to as the “ground truth” for the query. we can estimate $p(t/q)$ based solely on R_q . According to Bayes’ rule, $p(t/q)$ is

$$p(t/q) = \frac{p(q|t) \cdot p(t)}{p(q)} = \frac{p(q|t) \cdot p(t)}{\sum_{i \in \text{dates}(D)} p(q|i) \cdot p(i)}$$

$P(t)$ is the probability that day t contains a document, multinomially distributed over t , and $p(q)$ is the prior probability of finding a document relevant to q , and serves as a normalizing factor. Assume that we know the complete set of documents R that are relevant to q , we can directly estimate $p(q/t)$ using the distribution of the documents in R_q over time as follows:

$$p(q|t) = \frac{\text{count}(R_q, t)}{\text{count}(D, t)}$$

Using the Distribution of Matching Documents

To estimate $p(t/q)$ for a query q and time t in the absence of knowledge of the relevant documents for q and time t . We suggested using the top matching documents

for q and their relevance scores. Specifically, $p(t/q)$ is defined as a normalized and the weighted sum of the relevance scores of the top- k matching documents published at day t , we can seen as follows:

$$p(t/q) \approx \sum_{d \in D_{q;k}} p(t|d) \cdot \frac{p(d|q)}{\sum_{\hat{d} \in D_{q;k}} p(\hat{d}|q)}$$

where $D_{q;k}$ are the top- k documents from D for query q . To connect the temporal relevance value of a day t with days in the near past, we apply simple moving-average smoothing technique.

Using Binning

The retrieval model not only suggests the top- k matching documents as an approximation to the true relevant documents. The relevant document weights these documents based on their relevance scores. This direct dependency on the relevance scores for estimating the $p(t/q)$ values is somewhat problematic. We suggest a general framework to estimate $p(t/q)$ that addresses these

issues, so that it is less dependent on the underlying retrieval model by considering only the top-k matching documents without using their relevance scores directly. Algorithm describes our method to estimate the value

$p(q/t)$ of each time t for a given query q over a news archive D . We follow three basic steps:

Input: Query q , document collection D

Output: Time-based probability $p(q/t)$ for each time t

Step 1: Compute the query-frequency histogram for q using the publication time of the documents in D

Step 2: Partition the times into bins $b_0; \dots; b_j$ based on the histogram characteristics.

Step 3: Define the value $p(q/t)$ of each time t based on t 's bin, such that a time in b will have a higher value than a time in b if $i < j$.

As the first step of our approach, we produce a query-frequency histogram for a user query executed over a news archive by identifying all the documents in the archive that match the query. By using conjunctive Boolean semantics for matching is sufficient to draw an expressive histogram that approximates the real distribution of relevant time periods.

After generating the query-frequency histogram, we move to Step 2 of our approach and analyze the histogram to organize all times (days) into bins. We explore alternate binning techniques based on different underlying hypotheses on how to identify the important time intervals.

In Step 3 of our algorithm, we define the $p(q/t)$ values based on the assignment of times to bins $b_0; \dots; b_j$ from Step 2. Hence, we define the binning so that bin b_i should be associated with $p(q/t)$ values that are higher than those for bin b_j whenever $i < j$.

Using Word Tracking

To obtain the top-k matching documents for q as a first step, typically with $k \geq 500$. When performed on top of an unmodified search engine some of the processing needed to answer a search-engine query. Furthermore, query processing in a state-of-the-art search engine is often optimized to return the top-10 results and is not efficient for producing a larger set of results. To estimate $p(t/q)$ efficiently, we refine $p(q/t)$ one step further and assume independence between query terms, to get

$$p(t|q) \propto p(t) \cdot p(q|t) \approx p(t) \cdot \prod_{w \in q} p(w|t),$$

where $p(w/t)$ is the probability of generating query word w at day t .

4. INTEGRATING RELEVANCE IN SEARCH

We discussed the general family of time sensitive queries and claimed that “traditional” information retrieval engines do not take temporal relevance into account when answering these queries. The top-10 results that The New York Times search engine returns for the time-sensitive query. Sometimes queries issued over a news archive are after recent events or breaking news, as we discussed previously. Our approach processes a recently query by computing traditional topic relevance scores for each document, and then “boosting” the scores of the most recent documents. To estimate the relevance of a document d to a query q ; $p(d/q)$, the conditional probability that d is topically relevant to q is computed. In the original language models and in later modifications, the prior $p(d)$ is ignored since it is assumed to be uniform and constant for all documents.

In contrast to recency queries, for a general time-sensitive query, we do not know beforehand either the relevant time periods or the expected distribution of relevant documents over time. Language models are a state-of-the-art general approach for ranking documents in a collection according to their topic similarity with a query. The query likelihood model estimates the relevance of a document d to a query q by computing the conditional probability $p(d/q)$ that d is topically relevant to q . It is defined as the

$$p(d|q) \propto p(d) \cdot p(q|d).$$

To answer general time-sensitive queries, we want to identify not just the relevant documents for the query, but also the relevant time periods. We introduced a framework to complement the topical relevance of a document for a query with additional evidence. We build on this framework and on the idea of splitting a document d into a content component c_d as well as a temporal component t_d . We now describe a similar integration into the probabilistic relevance model, a leading state-of-the-art approach has been suggested. In defining PRM, we state the following principle: “To produce the optimal ranking of a set of documents as an answer for a Question the doc’s should be given ranking according the probability of belonging to the relevance class R of the query’s.” that are to be introduced the following general PRM framework:

$$p(R|d, q) \propto_q \log \frac{p(R|d, q)}{p(\bar{R}|d, q)} \propto_q \log \frac{p(d|R, q)}{p(d|\bar{R}, q)},$$

5. EXISTING SYSTEM

Search Engines uses an approach that might help it identify how relevant the results it displays to searchers which are generally likely those results are to show a variety of results when a searcher uses a query term that might cover a range of topics in future. Age old prior approaches used Human Reviewers being one option for checking on the relevancy of search results by manually screening the results for each query. Techniques for automatically checking the relevance and variety of search results are provided. Then a query is submitted to the search machine, which uses a finding algorithm to obtain search results based on the query. A set of the top n related terms for the query is identified. For each related term in the set of terms, then its relative frequency in relation to all terms in the set of terms is determined. If the term does not occur in any of the results, then a loss in variety proportional to the relative term frequency for the term has occurred. Otherwise, the relevance of the search results is calculated by comparing the proportion of results containing the term with the relative term frequency for a term. This process is repeated for all terms in the set of related terms to produce a total variety and relevance for the results. These processed and refined results are shown to the end user, for the query they’ve initiated. However these search algorithms consider relevancy as an important factor, they don’t initiate re-sorting based on time sensitivity to improve the results classification better. Unfortunately, ignoring or not fully exploiting the time dimension can be detrimental for a large family of queries for which we should consider not only the document topical relevance but the publication time of the documents. So a better system is required that can achieve that. We observe that, for an important class of queries over news archives that we call time-sensitive queries, topic similarity is not sufficient for ranking. For such queries, the publication time of the documents is important and should be considered in conjunction with the topic similarity to derive the final document ranking. Most current methods for searching over large archives of timed documents incorporate time in a relatively crude manner by only focusing on the published relevancy. Users can submit a keyword query, say [Madrid bombing], and restrict the results to articles written between March and April 2004, or alternatively sort the results on the publication date of the articles. Unfortunately, searchers do not always know the appropriate time intervals for their queries, and placing the burden on the users to explicitly handle time during querying is not desirable. Beyond asking for explicit user input, focus on handling recency queries, which are queries that are after recent events or breaking news amounting to high traffic in time line for that particular query. The Existing query model considers the following aspects. If the relevant time period for a time-sensitive query is unspecified, several query processing approaches are possible. One alternative is to automatically suggest, based on the query terms, relevant time ranges for the query and allow users to explicitly select appropriate time intervals. As an alternative that demands less input from the users, and which we follow in this paper, we can automate the previous procedure and prioritize results from periods that we automatically identify as relevant. We can then naturally define the relevance of a document as a combination of topic similarity and time relevance.

6. PROPOSED SYSTEM

For an important class of queries termed time-sensitive queries over frequently updated archives such as news archives, topic similarity alone is not sufficient for ranking. For such queries, the publication time of the documents is important and should be considered in conjunction with the

topic similarity to derive the final document ranking. For incorporating the time dimension, prior systems used publication date and time of the documents to locate time periods of interest. However, a document published at a later date (e.g., a review article, summarizing an event) may also be relevant; We propose to infer the temporal relevance of a document by analyzing its contents, and not by relying solely on its publication date thus increasing the relevancy of the results. So we propose to use Tf-idf, term frequency-inverse document frequency a numerical statistic method, that reflects how important a word is to a document in a collection or corpus. Recently being used as a weighting factor in information retrieval and text mining areas. The tf-idf value increases proportionally to the number of times a word appears in the document, but is offset by the frequency of the word in the corpus, which helps to control for the fact that some words are generally more common than others. Variations of the tf-idf weighting scheme is used by our search engine prototype as a central tool in scoring and ranking a document's relevance given a user query. tf-idf can be successfully used for stop-words filtering in various subject fields including text summarization and classification.

Our work demonstrates that integrating time in the retrieval task can improve the quality of the retrieval results, and motivates further research in the area.. It may also be relevant; an interesting direction for future research is to infer the temporal relevance of a document by analyzing its contents and not by relying solely on its publication date. We introduce time-based diversity in query results by grouping the results into clusters of relevant time ranges. Along the same lines, we are interested in integrating our retrieval techniques with algorithms for query reformulation. Overall, we believe that seamlessly integrating temporal information into web search for news articles or otherwise is a promising direction can significantly improve the web search experience. We design general framework to incorporate time into the retrieval task in a principled manner. These intervals are then used to adjust the document relevance scores by boosting the scores of documents published within the important intervals. Our system provides a web interface for searching the News blaster archive, an operational news archive and summarization system.

7. PERFORMANCE ANALYSIS

To answer the type of time-sensitive queries² over a news archive, we would like to use the temporal information implicitly available in the archive. Hence, we observe that time-sensitive queries are generally after documents from specific time periods.

This observation suggests that it is important to know the distribution of relevant documents over time for a given query. we use the publication time of the returned documents to generate the query frequency histogram over time. After generating the query-frequency histogram our approach is to explore alternate binning techniques based on different underlying hypotheses on how to identify the important time intervals. The “events” that these query likely targets last one or two days and, as a result, the query-frequency histogram of such a query will usually have sharp, thin “spikes” indicating these events. News events can last for longer than one or two days. An event appears in the query-frequency histogram, creating the shape of a “bump.” We compute the average daily query frequency in a window of x days into the past and x days into the future. We considered windows of fixed size around each day for binning. We identify continuous time intervals of variable length where the query frequency on each day is greater than the average query frequency per day in the entire collection. we define the binning so that bin b_j should be associated with $p(q/t)$.

8. RESULTS

We now report results for TQBLASTER, for BUMP-QL, BUMP-RM, SUM-QL, SUM-RM, QL-TOPIC, and RMTOPIC with the same $_$ values used for TQ351 and TQ401. We selected the best two baseline techniques and four time-sensitive techniques according to the TREC experiments, and excluded the other techniques to keep the amount of human annotations that we needed at manageable levels.

We have showed that considering time as an additional factor for ranking query results may be valuable for answering time-sensitive queries. Our results indicate that using temporal evidence derived from news archives often increases precision and reveals new relevant documents from important time intervals.

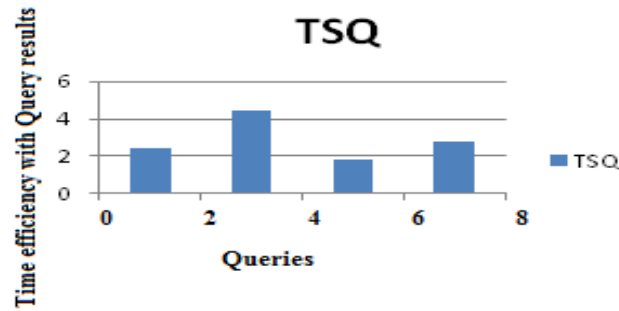


Figure 2. Performance results with time sensitive Queries

9. CONCLUSION

We presented a method for processing time-sensitive queries over a news archive, with techniques for identifying important time periods for a query. Our techniques improve the quality of search results as we presented an extensive experimental evaluation, including TREC as well as an archive of news articles. Our work demonstrates that integrating time in the retrieval task. Currently, we rely on the publication time of the documents to locate time periods of interest. Another promising research direction is to introduce time-based diversity in query results by grouping the results into clusters of relevant time ranges, enabling users to be aware of and interact with time information when examining the query results. Another interesting direction is to examine techniques that consider a time-sensitive definition of relevance at the document level. Of course, handling such a time-varying definition of relevance may require extensive rethinking of the existing ways of evaluating retrieval performance.

REFERENCES

- [1] Wisam Dakka, Luis Gravano, and Panagiotis G. Ipeirotis, "Answering General Time-Sensitive Queries", *Ieee Transactions On Knowledge And Data Engineering*, Vol. 24, NO. 2, February 2012
- [2] X. Li and W.B. Croft, "Time-Based Language Models," *Proc. 12th ACMConf. Information and Knowledge Management (CIKM '03)*, 2003
- [3] R. Jones and F. Diaz, "Temporal Profiles of Queries," *ACM Trans. Information Systems*, vol. 25, no. 3, article 14, 2007
- [4] S.E. Robertson, S. Walker, M. Hancock-Beaulieu, A. Gull, and M. Lau, "Okapi at TREC," *Proc. Fourth Text REtrieval Conf. (TREC-4)*, 1994.
- [5] S.E. Robertson, "Overview of the Okapi Projects," *J. Documentation*, vol. 53, no. 1, pp. 3-7, 1997.
- [6] K.S. Jones, S. Walker, and S.E. Robertson, "A Probabilistic Model of Information Retrieval: Development and Comparative Experiments Part 1," *Information Processing and Management*, vol. 36, no. 6, pp. 779-808, 2000.
- [7] K.S. Jones, S. Walker, and S.E. Robertson, "A Probabilistic Model of Information Retrieval: Development and Comparative Experiments Part 2," *Information Processing and Management*, vol. 36, no. 6, pp. 809-840, 2000.
- [8] I. Mani, J. Pustejovsky, and R. Gaizauskas, *The Language of Time: A Reader*. Oxford Univ. Press, 2005.
- [9] J.M. Ponte and W.B. Croft, "A Language Modeling Approach to Information Retrieval," *Proc. 21st Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '98)*, 1998.
- [10] F. Song and W.B. Croft, "A General Language Model for Information Retrieval," *Proc. Eighth ACM Conf. Information and Knowledge Management (CIKM '99)*, 1999.
- [11] N. Craswell, S.E. Robertson, H. Zaragoza, and M. Taylor, "Relevance Weighting for Query Independent Evidence," *Proc. 28th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '05)*, 2005.

- [12] S. Brin and L. Page, "The Anatomy of a Large-Scale Hypertextual Web Search Engine," Proc. Seventh Int'l World Wide Web Conf. (WWW '98), 1998.
- [13] V. Lavrenko and W.B. Croft, "Relevance-Based Language Models," Proc. 24th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '01), 2001.
- [14] S.E. Robertson, "The Probability Ranking Principle in IR," Readings in Information Retrieval, pp. 281-286, Morgan Kaufmann, 1997.
- [15] S.E. Robertson, S. Walker, and M. Hancock-Beaulieu, "Okapi at TREC-7: Automatic Ad Hoc, Filtering, VLC and Interactive Track," Proc. Seventh Text REtrieval Conf. (TREC-7), 1998.
- [16] N. Craswell, H. Zaragoza, and S.E. Robertson, "Microsoft Cambridge at TREC-14: Enterprise Track," Proc. 14th Text Retrieval Conf. (TREC-14), 2005.
- [17] K. McKeown, R. Barzilay, D. Evans, V. Hatzivassiloglou, J. Klavans, A. Nenkova, C. Sable, B. Schiffman, and S. Sigelman, "Tracking and Summarizing News on a Daily Basis with Columbia's Newsblaster," Proc. Second Int'l Conf. Human Language Technology (HLT '02), 2002.
- [18] R. Krovetz, "Viewing Morphology as an Inference Process," Proc. 16th Ann. Int'l ACM SIGIR Conf. Research and Development in Information Retrieval (SIGIR '93), 1993.