

On the Construction of Ideal Query for Patent Prior-art Search

M O
NICTA & ANU
Canberra, Australia
name.surname@nicta.com.au

W Z
NICTA & ANU
Canberra, Australia
name.surname@nicta.com.au

X Y
NICTA & ANU
Canberra, Australia
name.surname@nicta.com.au

ABSTRACT

Patent prior-art search aims to find all relevant patents which may invalidate the novelty of a patent application or at least have common parts with patent application and should be cited. Patent search has been the centre of attention in IR communities for years, however it has lower retrieval effectiveness compared to other IR applications. In this work, we focused on the causes of failure rather than solutions. We started with relevance feedback to get a golden standard, then we concentrated on heuristics correlate with our RF standard. Finally, we showed that features other than relevance feedback can not be helpful because they are a complex mixture of useful words and noisy words. Finally, we got a considerable improvement by user feedback with a minimum effort.

Categories and Subject Descriptors

H.3.3 [Information Search and Retrieval]: Query Formulation

Keywords

Patent search, Query Reformulation, Data Analysis

1. INTRODUCTION

A patent is a set of exclusive rights granted to an inventor to protect their invention for a limited period of time. An important requirement for a patent to be granted is that the invention, it describes, is novel which means there is no earlier patent, publication or public communication of a similar idea. To ensure the novelty of an invention, patent offices as well as other Intellectual Property (IP) service providers mainly perform a search called ‘prior art search’. The purpose of ‘prior art search’ is finding all relevant patents which may put the patent application at the risk of novelty invalidation or at least have common parts with patent application and should be cited [8] [14].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGIR '15, August 9-13, 2015, Santiago, Chile

Copyright 20XX ACM X-XXXXX-XX-X/XX/XX ...\$15.00.

Patent retrieval has three main characteristics which makes it difficult compared to other IR applications: the search starts with a query as long as a full patent application that helps users –usually patent examiners, inventors, or lawyers– avoid spending long hours to formulate a query; it is recall-oriented, where not missing relevant documents is more important than appearing relevant documents at top of the list; unlike the web application in which authors tend to highlight their work to be easily found through search engines, authors of the patents prefer to use a vague language to avoid the invalidation of their idea.

Many works has been conducted to improve the patent retrieval effectiveness so far. However, either the results showed quite small improvement or the proposed methods were complicated and computationally expensive. Overall, the works on patent search fall in five main categories: query reformulation(query expansion and query reduction), query term selection, query suggestions, using patent meta-data and images for retrieval [7], and Cross-Language Information Retrieval [10].

In this work, we mainly emphasized on the problem from the data analysis perspective rather than the solution, since the results, reported in the previous works, are lower than the performance for the other IR applications. We started with relevance feedback to find a golden standard for our analysis, then we examined possible recognized features to find a heuristic that correlate with our relevance feedback results. We avoided complex feature which are computationally expensive such as Pair-wise Term Proximity features [1]. Finally, we could double the ‘MAP’ with the minimum user effort.

2. BASELINE IR SYSTEM

We developed a Lucene-based¹ IR system with the possibility of using diverse generic IR models: TF-IDF, BM25, Language Models(Dirichlet smoothing, and Jelinek-Mercer smoothing) as our baseline system. We achieved the best baseline effectiveness using the ‘Description’ of the patent application as a query[15], and Language Model with Dirichlet smoothing as a retrieval model. We conducted our experiments on CLEF-IP²2010 data collection, with 2.6 million European patent documents and 1303 English topics(queries). On the collection side, we only indexed English subset of each section of a patent (title, abstract, claims, and description), and IPC³code in a separate field[8]. We also used

¹<http://lucene.apache.org/>

²<http://www.ifs.tuwien.ac.at/~clef-ip/>

³International Patent Classification

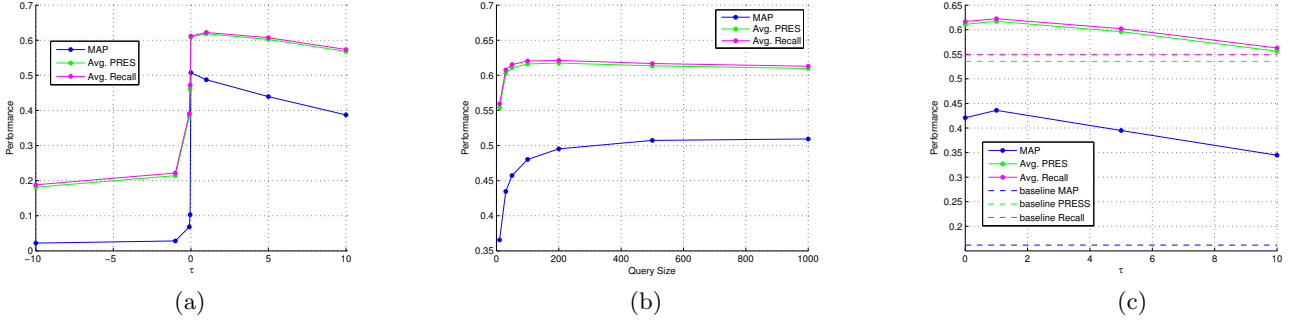


Figure 1: How score threshold(τ) and query size controls the performance. (a) Performance versus the score threshold. (b) Performance versus the query size. (c) System performance when we reduced the query by RF: $query = Q \cap (useful\ terms)$, where Q is the patent query and $useful\ terms = \{t | score_{RF}(t) > \tau\}$.

the patent classification assigned to the query topics to filter search results to match at least one of the query IPC codes[5]. Our experiments showed that using IPC filter is itself a source of error because about 19% of relevant patents in CLEF-IP data collection do not share any classification code with their query. However, for our analysis, we kept the filter on since it makes the matching process between the query and documents notably faster.

3. IDEAL QUERY

The main complain about patent search is insufficient match between the content of patent queries and relevant patents[6][8]. However, we have the intuition that there are sufficient terms in a patent query containing thousands words to be matched with the relevant patents. So, in this section, we focused on term analysis to figure out the main causes that the system fails in retrieving relevant documents at top of the result list.

We started our analysis using *relevance feedback*, in which the user gives feedback on the relevance of documents in an initial set of results to improve the final result set. We calculate a relevance feedback(RF) score for each term in top-100 retrieved documents as follows:

$$score_{RF}(t, Q) = Rel(t) - Irr(t) \quad (1)$$

$t \in \{\text{terms in top-100 retrieved documents}\}$

where $Rel(t)$ is the average term frequency in retrieved relevant patents and $Irr(t)$ is the average term frequency in retrieved irrelevant patents. We assumed that words with a positive score are *useful words* since they are more frequent in relevant patents, while words with negative score are *noisy words* as they appear more frequently in irrelevant patents.

We expected to see a higher performance for the queries which contain more *useful words*, but, surprisingly, we could not find any correlation between the performance and the percentage of *useful words* in the query.

3.1 Ideal Query Formulation

We hypothesized that a query, formulated by only the *useful terms*, is the best possible query we can make since they are all frequent in relevant patents but rare in irrelevant ones. We formulated the ideal query as follows:

$$Ideal\ query = \{t \in top - 100 | score_{RF}(t) > 0\} \quad (2)$$

Table 1 compares the baseline performance, where the query is the full patent application, with the performance of the ideal query. It can be seen that MAP jumps from 0.1618 to

Table 1: System performance for the baseline and ideal query.

	Pat.Query Weight:TF	Pat.Query Weight:1	Ideal Query Weight:Score(t)	Ideal Query Weight:1
PRES	0.5355	0.4268	0.6086	0.6087
MAP	0.1618	0.1181	0.4617	0.5075*
A. Recall	0.5491	0.4385	0.6129	0.6118

0.5075, which means the ideal query considerably performs better than the baseline. Hence, we used it as a golden standard.

3.2 Patent Query and Useful Terms

Our previous experiments led us to the hypothesis that a patent query contains sufficient words matched with the relevant patents. To prove our idea, we formulated a query by selecting only RF *useful terms* existing inside patent query as follows:

$$query = \{t | t \in \{Q \cap (useful\ terms)\}\} \quad (3)$$

The results were encouraging, as MAP was improved from 0.1618 to 0.44.

3.3 Analyse the Results

The main results related to ideal query formulation has summarized in Figure 1.

Fig. (1-c) explicitly shows that a patent query contains sufficient words to retrieve relevant patents at top of the list. We only need to keep the *useful terms* and prune out the *noisy words*.

We use a score threshold(τ) to formulate the RF query (we select the terms with $score(t) > \tau$).

Fig. (1-a) indicates two important facts. First, it shows that the performance decreases by increasing the threshold which means all words with a positive score are helpful for the performance. Second, the system is over-sensitive to the *noisy words*($\tau < 0$). Fig. (1-b) shows that formulating a query with up to 200 *useful words* helps performance whereas there is no significant improvement by adding more than 200 words.

4. STANDARD METHODS FOR QUERY REDUCTION

We achieved a notable improvement in performance using relevance feedback, however, it will cost to have a relevance feedback from our users who are professional patent examiners. For this reason, we considered the relevance feedback as a *golden standard* and focused to recognize the other accessible features that correlate with RF score. The most obvious accessible elements in our patent search engine are: patent query, top-100 patent documents retrieved in the first retrieval (features such as document frequent(DF) words and PRF words), and IPC code definition words.

4.0.1 Identify the Noisy Words

We showed that there are sufficient useful terms in a patent query(Fig. 1-c) to achieve good results but additional words which are considered noise are the main cause of low effectiveness. We hypothesized that the noisy words are frequent in top-100 retrieved documents and we can refine the query by removing document frequent words from the patent query. We calculated document frequent score for each term in top-100 retrieved patent documents as follows:

$$score_{DF}(t) = \frac{1}{100} \sum_{t \in \{Top-100\}} TF(t) \quad (4)$$

where $TF(t)$ is term frequency of each word. Fig. 2 shows

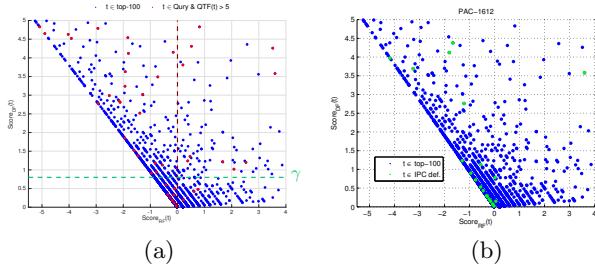


Figure 2: Anecdotal example: Scatter plot of $score_{RF}(t)$ versus $score_{DF}(t)$ for the words in top-100 retrieved documents. Each blue point is a vocabulary in top-100 retrieved document vocabulary set. (a) Points highlighted in red are query words with term frequency higher than 5 ($QTF(t) > 5$) and (b) points in green are IPC definition words.

a scatter plot of RF score and DF score of words in top-100 retrieved patents. A negative correlation shows that these two scores have opposite trends. Fig. 2-a is divided into four areas: (a) $score_{DF}(t) > \gamma$ and $score_{RF}(t) < 0$ (b) $score_{DF}(t) > \gamma$ and $score_{RF}(t) > 0$ (c) $score_{DF}(t) < \gamma$ and $score_{RF}(t) < 0$ (d) $score_{DF}(t) < \gamma$ and $score_{RF}(t) > 0$. If we remove the words with a DF score higher than a threshold ($score_{DF}(t) > \gamma$), we will remove noisy words of the area ‘a’ and keeping useful terms in area ‘d’ which is desirable. Whereas, the areas ‘b’ and ‘c’ are not desirable because many noisy words are kept in area ‘c’ while useful terms are removed from area ‘b’. To keep useful terms in area ‘b’, we added the second condition to remove DF words while keeping query terms with a term frequency higher than a threshold δ ($score_{DF}(t) < \gamma$ & $QTF(t) > \delta$). Even the best result, obtained with the $\gamma = 0.01$ and $\delta = 5$, could not beat the baseline. If we remove document frequent

words but keep query words with term frequency higher than 5 (highlighted with red in Fig. 2), we are still keeping sufficient amount of noisy words which are also frequent in the query to destroy the retrieval effectiveness.

To remove the noisy words which are both frequent in retrieved documents and query, we considered the third feature: IPC code definition. IPC code definition are short sentences giving a general description about the topic that each patent belongs to. In almost all queries IPC code definition words have a negative RF score but few with positive RF scores (Fig. 2-b). Our experiments showed that adding the third condition to remove frequent query word that is in IPC code definition words deteriorated the performance which means removing even just few useful words in patent query can be destructive.

4.0.2 Pseudo Relevance Feedback(PRF)

The main advantage of pseudo relevance feedback or a blind feedback is that it is an automated process without user interaction which assumes the top k ranked documents are relevant and the others are irrelevant. It has been found to improve performance in many applications, however it did not worked in patent prior-art search. The results for PRF query formulation and query term selection using PRF were below the baseline. In fact, we could not find any heuristic correlates between $score_{RF}(t)$ and $score_{PRF}(t)$. Fig. 3 explains the reason; it is an anecdotal example of a sample query with its abstract and a pair of PRF terms, with $score_{PRF}(t) > 10$, and RF score of each term. It can be seen that terms with high PRF score are considered noise since their RF score is negative. Fig. 3 combined with Fig. 1-a can justify that terms from PRF are not useful at all because they contain sufficient noisy words to destruct the retrieval effectiveness.

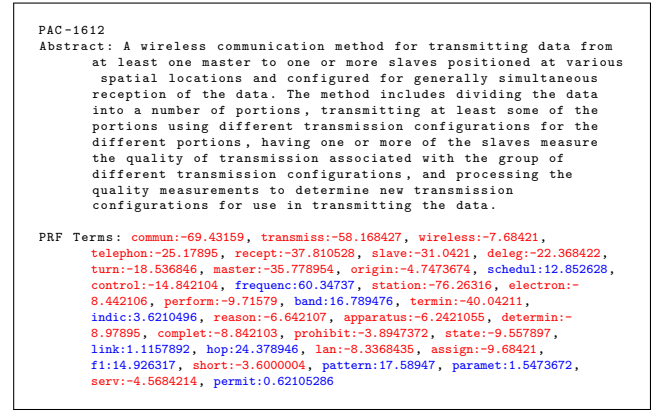


Figure 3: Anecdotal example: it shows the abstract and PRF term : $score_{RF}(PRF \text{ term})$ pair of a sample query. Useful terms are highlighted in blue and the noisy ones in red.

Surprisingly, there were no pattern or correlation between accessible features and RF score and we could not refine the patent query with just useful terms as what we could with relevance feedback.

5. IMPROVED BY MINIMUM USER EFFORT

All our attempts to improve the system effectiveness without accessing the relevance feedback were quite in vain because the features we recognized were tightly the combination of the useful words and noisy words and the system performance is too sensitive to the existence of a noisy word or the absence of the useful terms. So, we decided to involve the users but with the least effort. In this experiment, we selected the query words using merely the few first-ranked relevant patents. Table 2 shows that we can double the ‘MAP’ by only the first-ranked relevant document. We hy-

Table 2: System performance when only the first relevant patent used for query reduction. τ is RF score threshold, and k indicates the number of first relevant retrieved documents.

	$k = 1$ $\tau = 0$	$k = 1$ $\tau = 1$	$k = 3$ $\tau = 0$	$k = 3$ $\tau = 1$
PRES	0.4965	0.5016	0.5699	0.5727
MAP	0.3028	0.3040*	0.3879	0.3872
A. Recall	0.5040	0.5090	0.5757	0.5787

pothesised that recognising the first-ranked patent is easy for a patent examiner because we expected that it appears at top-5 in first retrieval. Fig. 4 confirms our intuition; The probability of finding the first-ranked relevant document at top-5 is acceptably high.

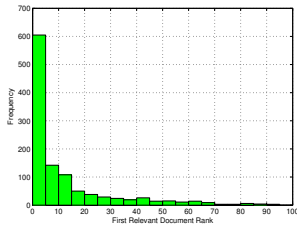


Figure 4: The distribution of the first relevant document rank over test queries which have TPs

6. RELATED WORK

Our work is different from pioneer studies on patent retrieval, as we closely looked into the problem rather than solutions to figure out the causes that generic IR models which are based on term matching process, do not work efficiently in patent domain. Magdy et al. [9] studied works on query expansion in patent retrieval and discussed that standard query expansion techniques are less effective, where the initial query is the full texts of query patents. Mahdabi et al. [12] used term proximity information to identify expansion terms. Ganguly et al. [2] adapted pseudo relevance feedback for query reduction by decomposing a patent application into constituent text segments and computing the Language Modelling (LM) similarities of each segment from the top ranked documents. The least similar segments to the pseudo-relevant documents removed from the query, hypothesizing it can increase the precision of retrieval. Kim et al. [3] provided diverse query suggestion using aspect identification from a patent query to increase the chance of retrieving relevant documents. Mahdabi et al. [11] used linked-based

structure of the citation graph together with IPC classification –the most useful patent meta-data– to improve the initial patent query.

7. CONCLUSIONS

In this paper, we looked at the patent prior-art search from a different perspective. While previous works proposed different solutions to improve retrieval effectiveness, we focused on term analysis of the patent query and top retrieved patents. After finding a golden standard from relevance feedback, we examined the most obvious features such as: document frequent words, query frequent words, IPC definition words, and pseudo relevance feedback that might correlate RF score for terms in top retrieved documents. We showed that these feature helps very little because they are a complicated mixture of useful terms and noisy words that can not be separated easily. Finally, we showed that we can double the ‘MAP’ with minimum user interaction. For future works, we plan to analyse more features which are independent from the relevance feedback but correlate with RF score. Inspired by some excellent works proposing query reduction and term selection techniques for the long non-patent queries[13][4], we are also going to apply them for patent retrieval.

8. ACKNOWLEDGMENTS

NICTA is funded by the Australian Government as represented by the Department of Broadband, Communications and the Digital Economy and the Australian Research Council through the ICT Centre of Excellence program.

9. REFERENCES

- [1] S. Bashir and A. Rauber. Improving retrievability of patents in prior-art search. In *Advances in Information Retrieval*, pages 457–470. Springer, 2010.
- [2] D. Ganguly, J. Leveling, W. Magdy, and G. J. Jones. Patent query reduction using pseudo relevance feedback. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 1953–1956. ACM, 2011.
- [3] Y. Kim and W. B. Croft. Diversifying query suggestions based on query documents. In *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pages 891–894. ACM, 2014.
- [4] G. Kumaran and V. R. Carvalho. Reducing long queries using query quality predictors. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 564–571. ACM, 2009.
- [5] P. Lopez and L. Romary. Patatras: Retrieval model combination and regression models for prior art search. In *Multilingual Information Access Evaluation I. Text Retrieval Experiments*, pages 430–437. Springer, 2010.
- [6] M. Lupu, A. Hanbury, et al. Patent retrieval. *Foundations and Trends in Information Retrieval*, 7(1):1–97, 2013.
- [7] M. Lupu, F. Piroi, and A. Hanbury. Evaluating flowchart recognition for patent retrieval. In *The Fifth International Workshop on Evaluating Information Access (EVIA)*, pages 37–44, 2013.

- [8] W. Magdy. *Toward higher effectiveness for recall-oriented information retrieval: A patent retrieval case study*. PhD thesis, Dublin City University, 2012.
- [9] W. Magdy and G. J. Jones. A study on query expansion methods for patent retrieval. In *Proceedings of the 4th workshop on Patent information retrieval*, pages 19–24. ACM, 2011.
- [10] W. Magdy and G. J. Jones. Studying machine translation technologies for large-data clir tasks: a patent prior-art search case study. *Information Retrieval*, 17(5-6):492–519, 2014.
- [11] P. Mahdabi and F. Crestani. Patent query formulation by synthesizing multiple sources of relevance evidence. *ACM Transactions on Information Systems (TOIS)*, 32(4):16, 2014.
- [12] P. Mahdabi, S. Gerani, J. X. Huang, and F. Crestani. Leveraging conceptual lexicon: query disambiguation using proximity information for patent retrieval. In *SIGIR*, 2013.
- [13] K. T. Maxwell and W. B. Croft. Compact query term selection using topically related text. In *Proceedings of the 36th international ACM SIGIR conference on Research and development in information retrieval*, pages 583–592. ACM, 2013.
- [14] F. Piroi, M. Lupu, and A. Hanbury. Overview of clef-ip 2013 lab. In *Information Access Evaluation. Multilinguality, Multimodality, and Visualization*, pages 232–249. Springer, 2013.
- [15] X. Xue and W. B. Croft. Transforming patents into prior-art queries. In *Proceedings of the 32nd international ACM SIGIR conference on Research and development in information retrieval*, pages 808–809. ACM, 2009.