

An Analysis of Query Reformulation Techniques for Precision Medicine

Maristella Agosti, Giorgio Maria Di Nunzio, Stefano Marchesin

Department of Information Engineering

University of Padua, Italy

{maristella.agosti,giorgiomaria.dinunzio,stefano.marchesin}@unipd.it

ABSTRACT

The Precision Medicine (PM) track at the Text REtrieval Conference (TREC) focuses on providing useful precision medicine-related information to clinicians treating cancer patients. The PM track gives the unique opportunity to evaluate medical IR systems using the same set of topics on two different collections: scientific literature and clinical trials. In the paper, we take advantage of this opportunity and we propose and evaluate state-of-the-art query expansion and reduction techniques to identify whether a particular approach can be helpful in both scientific literature and clinical trial retrieval. We present those approaches that are consistently effective in both TREC editions and we compare the results obtained with the best performing runs submitted to TREC PM 2017 and 2018.

CCS CONCEPTS

• **Information systems** → **Specialized information retrieval**; **Ontologies**; **Query reformulation**.

KEYWORDS

Medical IR; query reformulation; precision medicine

ACM Reference Format:

Maristella Agosti, Giorgio Maria Di Nunzio, Stefano Marchesin. 2019. An Analysis of Query Reformulation Techniques for Precision Medicine. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '19)*, July 21–25, 2019, Paris, France. ACM, New York, NY, USA, 4 pages. <https://doi.org/10.1145/3331184.3331289>

1 MOTIVATIONS

Medical Information Retrieval (IR) helps a wide variety of users to access and search medical information archives and data [4]. In [7, chapter 2], a classification of textual medical information is proposed: 1) Patient-specific information which applies to individual patients. This type of information can be structured, as in the case of an Electronic Health Record (EHR), or can be free narrative text. 2) Knowledge-based information that has been derived and organized from observational or experimental research. In the case of clinical research, the information is most commonly provided by books

and journals, but can take a wide variety of other forms, including computerized media. Therefore, the design of effective tools to access and search textual medical information requires, among other things, enhancing the query through expansion and/or rewriting techniques that leverage the information contained within knowledge resources. In this context, Sondhi et al. [12] identified some challenges arising from the differences between general retrieval and medical case-based retrieval. In particular, state-of-the-art retrieval methods, combined with selective query term weighing based on medical thesauri and physician feedback, improve performance significantly [3, 13].

In 2017 and 2018, the Precision Medicine (PM) [10] track¹ at the Text REtrieval Conference (TREC)² focused on an important use case in clinical decision support: providing useful precision medicine-related information to clinicians treating cancer patients. This track gives a unique opportunity to evaluate medical IR systems since the experimental collection is composed of a set of topics (synthetic cases created by precision oncologists) for two different collections that target two different tasks: 1) retrieving biomedical articles addressing relevant treatments for a given patient, and 2) retrieving clinical trials for which a patient – described in the information need – is eligible.

The objective of our study is to take advantage of this opportunity and evaluate several state-of-the-art query expansion and reduction techniques to examine whether a particular approach can be helpful in both scientific literature and clinical trials retrieval. Given the large number of participating research groups to this TREC track, we are able to compare the best experiments submitted to the PM track based on the results which were obtained applying our approach in the last two years. The experimental analysis shows that there are some common patterns in query reformulation that allow the retrieval system to achieve top performing results in both tasks.

The rest of the paper is organized as follows: Section 2 describes the approach used to evaluate different query reformulation techniques. Section 3 presents the experimental setup and compares the results obtained using our approach with the best performing runs from TREC PM 2017 and 2018. Finally, Section 4 reports some final remarks and concludes the paper.

2 APPROACH

The approach we propose for query expansion/reduction in a PM task comprises three steps, plus an additional fourth step required only for the retrieval of clinical trials. The steps are: (i) indexing, (ii) query reformulation, (iii) retrieval and (iv) filtering.

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '19, July 21–25, 2019, Paris, France

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-6172-9/19/07...\$15.00

<https://doi.org/10.1145/3331184.3331289>

¹<http://www.trec-cds.org/>

²<https://trec.nist.gov/>

Indexing Step. We create the following fields to index clinical trials collections: <docid>, <text>, <max_age>, <min_age> and <gender>. Fields <max_age>, <min_age> and <gender> contain information extracted from the eligibility section of clinical trials and are required for the filtering step. The <text> field contains the entire content of each clinical trial — and therefore also the information stored within the fields described above.

To index scientific literature collections, we create the following fields: <docid> and <text>. As for clinical trials, the <text> field contains the entire content of each target document.

Query Reformulation Step. The approach relies on two types of query reformulation techniques: query expansion and query reduction.

Query expansion: We perform a knowledge-based a priori query expansion. First, we rely on MetaMap [2], a state-of-the-art medical concept extractor, to extract from each query field all the Unified Medical Language System (UMLS)³ concepts belonging to the following semantic types⁴: Neoplastic Process (*neop*), Gene or Genome (*gngm*) and Cell or Molecular Dysfunction (*comd*). The *gngm* and *comd* semantic types are related to the query <gene> field, while *neop* is related to the <disease> field. For those collections where an additional <other> field is included — which considers other potential factors that may be relevant — MetaMap is used on <other> with no restriction on the semantic types, as its content does not consistently refer to any particular semantic type.

Second, for each extracted concept, we consider all its name variants contained into the following knowledge sources: National Cancer Institute⁵ (NCI), Medical Subject Headings⁶ (MeSH), SNOMED CT⁷ (SNOMEDCT) and UMLS Metathesaurus⁸ (MTH). All knowledge sources are manually curated and up-to-date.

The expanded queries consist in the union of the original terms with the set of name variants. For example, consider a query only containing the word “*melanoma*” — which is mapped to the UMLS concept C0025202. The set of name variants for the concept “*melanoma*” contains, among many others: cutaneous melanoma; malignant melanoma; malignant melanoma (disorder); etc. Therefore, the final expanded query is the union of the original term “*melanoma*” with all its name variants.

Additionally, we expand queries that do not mention any kind of blood cancer (e.g. “*lymphoma*” or “*leukemia*”) with the term *solid*. This expansion proved to be effective in [5] where the authors found that a large part of relevant clinical trials do not mention the exact disease. A more general term like *solid tumor* is preferable and more effective.

Query reduction: We reduce original queries by removing, whenever present, gene mutations from the <gene> field. To clarify, consider a topic where the <gene> field mentions “*BRAF (V600E)*”. After the reduction process, the <gene> field becomes “*BRAF*”. The reduction process aims to mitigate the over-specificity of topics, since the information contained in a topic is too specific compared to those contained in the target documents [8].

Additionally, we remove the <other> field from those collections that include it — since it contains additional factors that are not necessarily relevant, thus representing a potential source of noise in retrieving precise information for patients.⁹

Retrieval Step. We use BM25 [11] as retrieval model. Additionally, query terms obtained through query expansion are weighted lower than 1.0 to avoid introducing too much noise in the retrieval process [6].

Filtering Step. The eligibility section in clinical trials comprises, among others, three important demographic aspects that a patient needs to satisfy to be considered eligible for the trial, namely: minimum age, maximum age and gender; where minimum age and maximum age are the minimum and the maximum age, respectively, required for a patient to be considered eligible for the trial, while gender is the required gender.

Therefore, after the retrieval step, we filter out from the list of candidate trials those for which a patient is not eligible — i.e. his/her demographic data (age and gender) does not satisfy the three aforementioned eligibility criteria aforementioned. In those cases where part of the demographic data is not specified, a clinical trial is kept or discarded on the basis of the remaining demographic information. For instance, if the clinical trial does not specify a required minimum age, then it is kept or discarded based on its maximum age and gender required values.

3 SETUP AND EVALUATION

In this section, we describe the experimental collections and the setup used to conduct and evaluate our approach. Then, we compare the results obtained with our approach with those of the best performing systems from TREC PM 2017 and 2018. All these systems make use of external knowledge sources to enhance retrieval performance; moreover, most of them are complex multi-stage retrieval systems, like those proposed in [5, 8], while the approach we present is quite simple and straightforward — facilitating its reproducibility.¹⁰

Experimental Collections. Both tasks in TREC PM use the same set of topics, but with two different collections: scientific literature, clinical trials.

Topics consists of 30 and 50 synthetic cases created by precision oncologists in 2017 and 2018, respectively. In 2017, topics contain four key elements in a semi-structured format: (1) disease (e.g. a type of cancer), (2) genetic variants (primarily present in tumors), (3) demographic information (e.g. age, gender), and (4) other factors (which could impact certain treatment options). In 2018, topics contain three of the four key elements used in 2017: (1) disease, (2) genetic variants, and (3) demographic information.

Scientific Literature consists of a set of 26,759,399 MEDLINE¹¹ abstracts, plus two additional sets of abstracts: (i) 37,007 abstracts from recent proceedings of the American Society of Clinical Oncology (ASCO), and (ii) 33,018 abstracts from recent proceedings of the American Association for Cancer Research (AACR). These

³<https://www.nlm.nih.gov/research/umls/>

⁴<https://metamap.nlm.nih.gov/SemanticTypesAndGroups.shtml>

⁵<https://www.cancer.gov/>

⁶<https://www.ncbi.nlm.nih.gov/mesh/>

⁷<https://www.snomed.org/>

⁸https://www.nlm.nih.gov/research/umls/knowledge_sources/metathesaurus/

⁹In a personal communication with the organizers of the track, we have been informed that it was difficult to convince the oncologists why the *other* field was even necessary.

¹⁰Source code available at: https://github.com/stefano-marchesin/TREC_PM_qreforms

¹¹<https://www.nlm.nih.gov/bsd/pmresources.html>

additional datasets were added to increase the set of potentially relevant treatment information. In fact, precision medicine is a fast-moving field where keeping up-to-date with the latest literature can be challenging due to both the volume and velocity of scientific advances. Therefore, when treating patients, it may be helpful to present the most relevant scientific articles for an individual patient. Relevant literature articles can guide precision oncologists to the best-known treatment options for the patient's condition.

Clinical Trials consists of a total of 241,006 clinical trial descriptions, derived from ClinicalTrials.gov¹² — a repository of clinical trials in the U.S. and abroad. When none of the available treatments are effective on oncology patients, the common recourse is to determine if any potential treatments are undergoing evaluation in a clinical trial. Therefore, it would be helpful to automatically identify the most relevant clinical trials for an individual patient. Precision oncology trials typically use a certain treatment for a certain disease with a specific genetic variant (or set of variants). Such trials can have complex inclusion and/or exclusion criteria that are challenging to match with automated systems.

Experimental Setup. We use Whoosh,¹³ a pure Python search engine library, for indexing, retrieval and filtering steps. For BM25, we keep the default values $k_1 = 1.2$ and $b = 0.75$ provided by Whoosh — as we found them to be a good combination [1]. For query expansion, we rely on MetaMap to extract and disambiguate concepts from UMLS. We summarize the procedure used for each experiment below.

Indexing

- Index clinical trials using the following created fields: <docid>, <text>, <max_age>, <min_age> and <gender>;
- Index scientific abstracts using the following created fields: <docid> and <text>.

Query reformulation

- Use MetaMap to extract from each query field the UMLS concepts restricted to the following semantic types: *neop* for <disease>, *gngm*/*cmd* for <gene> and *all* for <other>;
- Extract from the concepts all name variants belonging to NCI, MeSH, SNOMED CT and MTH knowledge sources;
- Expand (or not) topics that do not mention “lymphoma” or “leukemia” with the term *solid*;
- Reduce (or not) queries by removing, whenever present, gene mutations from the <gene> field;
- Remove (or not) the <other> field.

Retrieval

- Adopt any combination of the reformulation strategies;
- Weigh expanded terms with a value $k \in \{0, 0.1, 0.2, \dots, 1\}$;
- Perform a search using expanded queries with BM25.

Filtering

- Filter out clinical trials for which the patient is not eligible.

Evaluation Measures. We use the official measures adopted in the TREC PM track: inferred nDCG (infNDCG), R-precision (Rprec) and Precision at rank 10 (P₁₀). Precision at rank 5 and at rank 10 were used only for the Clinical Trials task 2017 and are not

reported in this work for space reasons. The inferred nDCG was not computed for the task Clinical Trials 2017 since the sampled relevance judgments are not available.

Comparison. In Table 1, we report the results of our experiments (upper part) and compare them with the top performing participants at TREC 2017 and 2018 (lower part). Given the large number of experiments, we decided to present the top 5 runs ordered by P₁₀ for each year and for each task. Each line shows a particular combination (*yes* or *no* values) of semantic types (*neop*, *cmd*, *gngm*), usage and expansion of <other> field (*oth*, *oth_exp*), query reduction (*orig*), and expansion using weighted *solid* (tumor) keyword. We use the symbol ‘.’ to indicate that the features *oth*, *oth_exp* are not applicable for year 2018 due to the absence of the <other> field in 2018 topics. We report the results for both Scientific Literature (*sl*) and Clinical Trials (*ct*) tasks. We highlight in bold the top 3 scores for each measure, and we use the symbols † and ‡ to indicate two combinations that performed well in both 2017 and 2018. For the TREC PM participants, we select those participants who submitted runs in both years and reached the top 10 performing runs in at least two measures [9, 10]. The results reported in the lower part of Table 1 indicate the best score obtained by a particular run for a specific measure; the best results of a participant are often related to different runs. The symbol ‘-’ means that the measure is not available, while ‘<’ indicates that none of the runs submitted by the participant achieved the top 10 performing runs. For comparison, we add for each measure the lowest score required to enter the top 10 TREC results list, and the score obtained by the best combination of our approach — indicated by the line number — as if we were participants of these tracks.

In 2018, there is a clear distinction in terms of performances among the combinations that achieve the best results for the *sl* and the *ct* tasks. For the *sl* task, considering the semantic type *neop* expansion without using the umbrella term *solid* provides the best performances for all the measures considered. On the other hand, two of the best three runs for the *ct* task (line 5 and 9), use no semantic type expansion, but rely on the *solid* (tumor) expansion with weight 0.1.

In 2017, the situation is completely different. Lines 12 and 13 show two combinations that are in the top 3 performing runs for both *sl* and *ct*. These two runs use query reduction and a weighted 0.1 *solid* (tumor) expansion. The use of a weighted 0.1 *solid* expansion as well as a reduced query (*orig* = *n*) seems to improve performances consistently for all measures in 2017. The semantic type *gngm* seems more effective than *neop*, while *cmd* does not seem to have any positive effect at all.

Another element that shows how difficult these two tasks are is the fact that top performing systems in 2017 do not achieve the same results in 2018. Our study therefore helps researchers to select (or remove) semantic types to build strong baselines for both tasks.

4 CONCLUSIONS AND FINAL REMARKS

In this paper, we proposed and evaluated several state-of-the-art query expansion and reduction techniques for scientific literature and clinical trials retrieval. The experimental analysis showed that no clear pattern emerges for both tasks. In general, a query expansion approach using a selected set of semantic types helps the

¹²<https://clinicaltrials.gov/>

¹³<https://whoosh.readthedocs.io/en/latest/intro.html>

Table 1: Results for the TREC PM tasks 2017 and 2018. Details are reported in Section 3.

976