

DCU@TRECMed 2012: Using ad-hoc baselines for domain-specific retrieval

Johannes Leveling, Lorraine Goeriot, Liadh Kelly, Gareth J. F. Jones

Centre for Next Generation Localisation (CNGL), School of Computing,
Dublin City University, Dublin 9, Ireland
`firstname.lastname@computing.dcu.ie`

Abstract. This paper describes the first participation of DCU in the TREC Medical Records Track (TRECMed) 2012. We performed initial experiments on the the 2011 TRECMed data based on the BM25 retrieval model. Surprisingly, we found that the standard BM25 model with default parameters performs comparable to the best automatic runs submitted to TRECMed 2011 and our experiments would have ranked among the top four out of 29 participating groups. We expected that some form of domain adaptation would increase performance. However, results on the 2011 data proved otherwise: query expansion decreased performance, and filtering and reranking by term proximity also decreased performance slightly. We submitted four runs based on the BM25 retrieval model to TRECMed 2012 using standard BM25, standard query expansion, result filtering, and concept-based query expansion. Official results for 2012 confirm that domain-specific knowledge, as applied by us, does not increase performance compared to the BM25 baseline.

1 Introduction

This paper describes the first participation of DCU in the TREC Medical Records Track (TRECMed 2012). TRECMed is an instance of domain-specific information retrieval (IR) and ran for the first time in 2011, with 29 participating groups. The second, and final, TRECMed track ran in 2012, with 24 participating groups. The search task within the TRECMed is an ad hoc search task that models the clinical task of finding cohorts for comparative effectiveness research, based on a document collection of de-identified clinical reports.

A review of the 2011 participants' approaches shows that the most successful automatic approaches include:

- Information extraction on the corpus: applying natural language processing (NLP) techniques (e.g. chunking, PoS tagging, lemmatization) [1–3]; linking with additional concept bases, or medical ontologies [1–5]; expanding ICD9 codes (International Statistical Classification of Diseases and Related Health Problems, version 9) for the patient's admission or discharge status [1, 5]; treating negation (e.g. negative test results or symptoms) [1–5];
- Query expansion based on external knowledge (e.g. medical web sites, knowledge bases or Wikipedia) [2, 3];

- Result filtering based on extracted patient features such as ethnicity, age, and gender [1, 4, 5];
- A popular IR framework used in the 2011 campaign is the Lucene toolkit (and its default retrieval model) [1–5].

Similar to TREC Med 2011, in 2012 the most successful automatic approaches included information extraction, query expansion and result filtering approaches [6–10]. Query expansion using pseudo-relevance feedback was also used in the 2012 task [8]. However, while in 2011 the default retrieval model in the Lucene toolkit (a vector space model style retrieval algorithm) was used by the top performing participants, in 2012 more sophisticated retrieval approaches were explored by some participants. For example, [7] investigated using the query likelihood language model and the Markov random field model for retrieval, which coupled with information extraction and filtering techniques yielded the best automatic results in TREC Med 2012.

Our goal for the participation in TREC Med 2012 was to investigate query processing and different expansion techniques while trying to establish a good baseline for this task.

2 Related Work

We view medical record retrieval as an instance of domain-specific IR. There have been several evaluation tasks in IR evaluation campaigns such as TREC, NTCIR, and CLEF which focus on domain-specific IR, e.g. TREC-Chem¹ [11], patent retrieval² [12], or geographic IR³ [13].

Domain adaptation for IR has not yet proven to be consistently successful. In GIRT, the domain-specific IR task at CLEF, few participants used meta-data such as additional document fields containing subject terms or a domain-specific thesaurus [14, 15], as standard IR models yielded a high performance. Similarly, one important result from evaluation of geographic IR (GIR) is that adding large gazetteers with geographic knowledge decreases performance. For GIR, simple text-based retrieval (with a bag-of-words approach) turned out to be a very strong baseline [16].

Armstrong et al. [17] analyzed several years of experiments on TREC data and found that very few results reported in the literature outperform strong baseline experiments. Most experimenters claim a significant improvement, but over a weak baseline and not over the best results on the same data. Thus, improvements on the data do not add up, as it becomes more difficult to improve on good results.

¹ http://wiki.ir-facility.org/index.php/TREC_Chemistry_Track

² <http://www.cl.cs.titech.ac.jp/~fujii/ntc8patmt/rs.html>

³ <http://www.linguateca.pt/GikiCLEF/>

3 System Description

The objective of our participation in TREC Med 2012 was to establish a baseline BM25 system and compare different query expansion techniques. Our system employs approaches that were described as successful by last year's participants and comprises simple preprocessing and analysis steps:

- The document terms from an initial indexing run were manually examined and a list of misspelled and run-together words was compiled. This list was used to correct terms in the final indexing stage.
- A single text index for fields was employed, formed from the report text and the textual description of the ICD9 fields.
- All report documents are indexed separately, i.e. retrieved results have to be mapped to patient visits. The system then returns the document with the maximum score to map reports to visits. This approach was found to best the best way to map retrieved documents to visits [3].
- Retrieval and query expansion are based on the BM25 model [18].
- An additional filtering step filters results by applying constraints from the query pertaining to the patient's age, ethnicity, gender, or admission status to the result set. Similar approaches have been investigated for TREC Med 2011 [1, 2, 5].

3.1 Document preprocessing

All report files (which pertain to a patient's visit) were indexed as separate documents. ICD9 codes were mapped to a description of the code, usually a short phrase/sentence. For instance, the ICD9 code *253.5* corresponds to the disease *Diabetes insipidus*. The code descriptions were then stored as additional document fields ICD9_DIS_DIAGNOSIS_TEXT and ICD9_ADM_DIAGNOSIS_TEXT. The fields REPORT_TEXT, TYPE, SUBTYPE, ICD9_DIS_DIAGNOSIS_TEXT, and ICD9_ADM_DIAGNOSIS_TEXT were used to create a single index for the body of text.

3.2 Spelling correction

Spelling errors may have a detrimental impact on the system's performance. We manually corrected terms by examining all index terms from an initial run. Corrections were only added to the list when the correction was unambiguous. The correction was not performed by a medical expert, so many incorrect technical terms may have been missed. Even so, we compiled a list of 9533 spelling errors from the medical documents, which was added to a list of 4192 frequent spelling errors compiled from Wikipedia. During indexing, misspelled words are replaced with their corrections from this list. As an example, we found eight misspellings for the word *admission*: *admision*, *admision*, *admissin*, *admissoin*, *admisson*, *admission*, *admsission*, *dmission*.

3.3 Retrieval

Most implementations of BM25 for Lucene approximate document length as the number of characters, approximate field length by the maximum field length (for BM25F), or store the length information with a loss of precision (e.g. [19]). This can result in lower performance and/or different optimal model parameters. We employ our own BM25 implementation for Lucene which follows the original BM25 description [18] closely. Our system employed Lucene’s standard tokenization and a standard stopword list containing 33 stopwords.

3.4 Query expansion

We applied two approaches to query expansion: the standard approach described by Robertson et al. [18] and concept expansion.

- Query expansion: In the default query expansion approach, terms from the top ranked documents are ranked by a term selection value [20] and the top R terms are added to the query. For our experiments, 10 terms were extracted from the top 10 documents.
- Concept expansion: We annotated the TRECMed queries with concepts from the UMLS (Unified Medical Language System) thesaurus⁴, using the MetaMap system [21]. For each phrase, the system gives a ranked list of potential mapping concepts (called *Meta Candidates*) and one or several (in case of equal scores) mapped concepts (called *Meta Mapping*). We used the *Meta Mapping* concept list and their short description to extend the query, for example: *Patients with complicated GERD who receive endoscopy* will be extended with *Gastroesophageal reflux disease, Clinic / Center - Endoscopy*

3.5 Result filtering

Initial retrieval results were filtered with respect to constraints given in the query regarding the age, gender, ethnicity, and admission status of a patient. Sentences containing the anonymized age information of a patient (*** age*) were extracted from the document collection and manually annotated to obtain annotation patterns. Roughly 500 patterns were extracted. For example, the word sequence “*is an ** age [in 80s] yr old wm admitted*” allows to infer that an 80-89 year old white male patient was admitted. The longest match between sentences in a document and the patterns was then used to augment the document’s meta-data and overwrite the default value (e.g. *unknown*) for the *age*, *gender*, *ethnicity*, and *admission status* features. For the filtering step, only documents with exactly matching values were kept, while allowing the value *unknown* to match with any value. Table 1(a), 1(b), 1(c), and 1(d) show the distribution of values in the annotated document collection.

⁴ <http://www.nlm.nih.gov/research/umls/>

Table 1. Distribution for extracted patient features.

| (a) <i>age</i> | | (b) <i>gender</i> | | (c) <i>ethnicity</i> | | (d) <i>admission status</i> | |
|----------------|--------|-------------------|--------|----------------------|--------|-----------------------------|--------|
| age | count | gender | count | ethnicity | count | admission status | count |
| 0-12 | 37 | female | 16.824 | asian | 18 | admitted | 12.369 |
| 13-20 | 529 | male | 14.592 | black | 1.400 | not admitted | 8.222 |
| 20-29 | 2.684 | unknown | 64.285 | hispanic | 4 | unknown | 75.110 |
| 30-39 | 2.675 | Σ | 95.701 | white | 5.885 | Σ | 95.701 |
| 40-49 | 5.385 | | | unknown | 88.394 | | |
| 50-59 | 6.611 | | | Σ | 95.701 | | |
| 60-69 | 6.561 | | | | | | |
| 70-79 | 6.661 | | | | | | |
| 80-89 | 6.418 | | | | | | |
| 90+ | 788 | | | | | | |
| unknown | 57.352 | | | | | | |
| Σ | 95.701 | | | | | | |

4 Experiments

4.1 Experiments on 2011 Data

We performed initial experiments on the the 2011 TRECMed data based on the previously described setup. We tested four runs on the 2011 data, using

- i) standard BM25 retrieval,
- ii) i) + standard query expansion (QE),
- iii) ii) + result filtering, and
- iv) ii) + concept-based query expansion (CE).

As an official evaluation measure, bpref was used in the official TRECMed 2011 runs, due to problems associated with calculating inferred measures. Results of our runs are shown in Table 2. Surprisingly, we found that our baseline experiment, applying the standard BM25 model with default parameters, performs comparable to the best automatic runs submitted to TRECMed 2011. It would have ranked among the top five out of 29 participating groups (0.4052 MAP, 0.5082 bpref, 0.6 P@10). Evaluation results for the best automatic runs range from 0.552-0.494 bpref, 0.656-0.568 P@10, and 0.440-0.401 Rprec for the top participants in TRECMed 2011 [22].

Domain-specific IR typically requires adaptation of at least one component of a search system to the domain, e.g. by including domain-knowledge from ontologies or modifying the retrieval model. We expected that some form of domain adaptation would increase performance compared to the BM25 retrieval baseline. However, results on the 2011 data did not confirm this: standard query expansion decreases MAP and bpref, but slightly increased precision at early ranks; concept-based query expansion decreased performance in general, and filtering and reranking results also decreased performance.

For comparison, we also report the median numbers of all TREC Med submissions (where available to use). The median could serve as a baseline result in itself, but it typically does not measure the effect of domain adaptation in isolation.

Our observations are in contrast to official results reported for TREC Med 2011, where the domain adaptation techniques were reported to increase performance; these effects could thus be attributed to weaker baselines.

Table 2. Results on 2011 topics.

| Run Description | MAP | bpref | P@10 | infAP | infNDCG | Rprec |
|---------------------|--------|--------|--------|--------|---------|--------|
| TREC median | - | 0.4115 | 0.4764 | - | - | 0.3087 |
| i) BM25 | 0.4052 | 0.5082 | 0.6000 | 0.3228 | 0.5802 | 0.4112 |
| ii) BM25+QE | 0.3249 | 0.4867 | 0.4853 | 0.3044 | 0.5717 | 0.3566 |
| iii) BM25+QE+filter | 0.3229 | 0.4857 | 0.4824 | 0.3037 | 0.5703 | 0.3549 |
| iv) BM25+CE+filter | 0.3425 | 0.5116 | 0.4882 | 0.3016 | 0.5411 | 0.3705 |

4.2 Experiments on 2012 Data

Results for our four official submitted runs are shown in Table 3. TREC Med 2012 used inferred evaluation measures. For comparison with the results we obtained on the 2011 collection, we also show the bpref results. We expected that, as for the 2011 data, concept-based query expansion and result filtering will decrease performance (bpref and MAP) significantly, compared to the simple BM25 baseline. We observe that performance in general is much lower compared to results on 2011 data, which may be due to more difficult topics. However, comparing our own experiments, the expected decreases is not as high as on 2011 data, e.g. run iii) vs. run i). Further, we did not rank amongst the top performing groups in 2012. The best automatic runs in 2012 ranged from 0.578-0.509 infNDCG, 0.286-0.231 infAP, and 0.592-0.553 P@10 for the top five participants in TREC Med 2012 [23].

Table 3. Results on 2012 topics.

| Run Description | MAP | bpref | P@10 | infAP | infNDCG | Rprec |
|---------------------|--------|--------|--------|--------|---------|--------|
| TREC median | - | - | 0.4702 | 0.1689 | 0.4244 | 0.2961 |
| i) BM25 | 0.2930 | 0.3462 | 0.4638 | 0.2069 | 0.4043 | 0.3135 |
| ii) BM25+QE | 0.2562 | 0.3163 | 0.4213 | 0.1784 | 0.3766 | 0.2912 |
| iii) BM25+QE+filter | 0.2734 | 0.3331 | 0.4553 | 0.1879 | 0.3947 | 0.3045 |
| iv) BM25+CE+filter | 0.2552 | 0.3152 | 0.4191 | 0.1784 | 0.3745 | 0.2953 |

5 Discussion

We see several possible explanations for the results we observed:

1. Our approach to include medical knowledge performs worse than the approaches of other participants because our annotation is less accurate. Using additional domain information with low accuracy degrades the performance. However, this would not explain why BM25 with default settings performs exceptionally well. In fact, the accuracy for the extracted patient features must be comparatively high, as the annotation patterns were manually extracted. Also, the distribution of the patient features is similar to that reported by other participants in 2011.
2. The BM25 retrieval model is superior to Lucene’s internal ranking scheme, which is a variant of tf-idf with support for boosting terms and documents. BM25 can still be considered a strong baseline, even for domain-specific IR and twenty years after its introduction. Lucene (with its standard ranking model) was used by many of the top performing groups in TRECMed 2011, but less so in 2012.
3. The query expansion methods were not optimized for TRECMed and thus, showed no improvement over the baseline. We performed additional experiments for the standard feedback approach (not reported in this paper), where we compared the performance of extracting feedback terms from the top ranked visits to top ranked reports. Results were found to be not significantly different from the results reported in this paper. The argument above is at least partially true for the concept-based query expansion approaches, where we did not explicitly control term weighting (e.g. by downweighting expansion terms) and did not limit the number of feedback terms.

In summary, BM25 proves to be a strong baseline for 2011 data. For 2012 data, performance for infAP and Rprec achieved with the standard BM25 model still exceeds the median values.

6 Conclusion

Including domain-specific adaptation results in more complex indexing and retrieval workflows, but intuitively, adaptation should result in a significant performance improvement over the standard retrieval baseline. For ad-hoc IR, Armstrong et al. [17] pointed out that comparing against a weak baseline allows observing a significant performance increase that cannot be replicated against a strong baseline. We argue that for a domain-specific task, strong baselines are needed even more to isolate domain-adaptation issues and make their effects observable. Strong generic baselines can be derived from open-domain retrieval (i.e. ad-hoc IR). Weak baselines are not adequate and invalidate conclusions on the effect of domain adaptation, because an improvement over a weak baseline is harder to reproduce over a stronger baseline.

We would like to propose for domain-specific information retrieval tasks in general, that additional meta-data or annotations are made available by participants in evaluation tasks as stand-off annotations for the document collection so that participating groups can perform experiments on the same meta-data (e.g. document ID and extracted patient features). This would lower the entry-level for new participants and make results between participants more comparable, as the quality of generating additional meta-data would be separate from the quality of using it. Furthermore, instead of computing artificial baselines, task organizers should provide baselines based on results obtained with state-of-the-art retrieval models, but without domain adaptation.

Acknowledgment

This research is supported by the Science Foundation Ireland (Grant 07/CE/I1142) as part of the Centre for Next Generation Localisation (CNGL) project at DCU and by funding from the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement n°257528 (KHRESMOI).

References

1. King, B., Wang, L., Provalov, I., Zhou, J.: Cengage Learning at TREC 2011 medical track. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
2. Goodwin, T., Rink, B., Roberts, K., Harabagiu, S.M.: Cohort shepherd: discovering cohort traits from hospital visits. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
3. Schuemie, M., Trieschnigg, D., Meij, E.: DutchHatTrick: Semantic query modeling, ConText, session detection, and match score maximization. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
4. Demner-Fushman, D., Abhyankar, S., Jimeno-Yepes, A., Loane, R., Rance, B., Lang, F., Ide, N., Apostolova, E., Aronson, A.R.: A knowledge-based approach to medical records retrieval. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
5. Gurulingappa, H., Müller, B., Hofmann-Apitius, M., Fluck, J.: A semantic platform for information retrieval from e-health records. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
6. Tinsley, B., Thomas, A., McCarthy, J., Lazarus, M.: Atigeo at TREC 2012 medical records track: ICD-9 code description injection to enhance electronic medical record search accuracy. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)
7. Rabiou, A.B., Chandar, P., Kumar, N., Rao, A., Zhu, D., Carterette, B.: UDel (Carterette) at TREC 2012. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)
8. Qi, Y., Laquerre, P.F.: Retrieving medical records: NEC Labs America at TREC 2012 medical record track. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)
9. P., M.C., Wang, Y., Fang, H.: Exploiting domain thesaurus for medical record retrieval. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)
10. Limsopatham, N., McCreddie, R., Albakour, M.D., Macdonald, C., Santos, R., Ounis, I.: University of Glasgow at TREC 2012: Experiments with Terrier in medical records, microblog, and web tracks. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)

11. Lupu, M., Piroi, F., Huang, X., Zhu, J., Tait, J.: Overview of the TREC 2009 Chemical IR Track. In Voorhees, E.M., Buckland, L.P., eds.: Proceedings of The Eighteenth Text REtrieval Conference, TREC 2009, Gaithersburg, Maryland, USA, November 17-20, 2009. Volume Special Publication 500-278., NIST (2009)
12. Fujii, A., Iwayama, M., Kando, N.: Overview of the patent retrieval task at the NTCIR-6 workshop. In: Proceedings of NTCIR-6 Workshop Meeting. (2007) 359–365
13. Santos, D., Cardoso, N., Carvalho, P., Dornescu, I., Hartrumpf, S., Leveling, J., Skalban, Y.: GikiP at GeoCLEF 2008: Joining GIR and QA forces for querying Wikipedia. In: Evaluating Systems for Multilingual and Multimodal Information Access, 9th Workshop of the Cross-Language Evaluation Forum, CLEF 2008, Aarhus, Denmark, September 17-19, 2008, Revised Selected Papers. Volume 5706 of LNCS., Springer (2009) 894–905
14. Kluck, M., Stempfhuber, M.: Domain-specific track CLEF 2005: Overview of results and approaches, remarks on the assessment analysis. In: Accessing Multilingual Information Repositories, 6th Workshop of the Cross-Language Evaluation Forum, CLEF 2005, Vienna, Austria, 21-23 September, 2005, Revised Selected Papers. Volume 4022 of LNCS., Springer (2006) 212–221
15. Petras, V., Baerisch, S., Stempfhuber, M.: The domain-specific track at CLEF 2007. In: Advances in Multilingual and Multimodal Information Retrieval, 8th Workshop of the Cross-Language Evaluation Forum, CLEF 2007, Budapest, Hungary, September 19-21, 2007, Revised Selected Papers. Volume 5152 of LNCS., Springer (2008) 160–173
16. Mandl, T., Carvalho, P., Nunzio, G.M.D., Gey, F.C., Larson, R.R., Santos, D., Womser-Hacker, C.: GeoCLEF 2008: The CLEF 2008 cross-language geographic information retrieval track overview. In: Evaluating Systems for Multilingual and Multimodal Information Access, 9th Workshop of the Cross-Language Evaluation Forum, CLEF 2008, Aarhus, Denmark, September 17-19, 2008, Revised Selected Papers. Volume 5706 of LNCS., Springer (2009) 808–821
17. Armstrong, T.G., Moffat, A., Webber, W., Zobel, J.: Improvements that don't add up: ad-hoc retrieval results since 1998. In: Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009, ACM (2009) 601–610
18. Robertson, S.E., Walker, S., Jones, S., Hancock-Beaulieu, M.M., Gatford, M.: Okapi at TREC-3. In Harman, D.K., ed.: Overview of the Third Text Retrieval Conference (TREC-3), Gaithersburg, MD, USA, NIST (1995) 109–126
19. Pérez-Iglesias, J., Pérez-Agüera, J.R., Fresno, V., Feinstein, Y.Z.: Integrating the Probabilistic Models BM25/BM25F into Lucene. CoRR (2009)
20. Walker, S., Robertson, S.E., Boughanem, M., Jones, G.J.F., Jones, K.S.: Okapi at TREC-6 Automatic ad hoc, VLC, routing, filtering and QSDR. In: TREC-6. (1997) 125–136
21. Aronson, A.R., Lang, F.M.: An overview of MetaMap: historical perspective and recent advances. *Journal of the American Medical Informatics Association* **17** (2010) 229–236
22. Voorhees, E.M., Tong, R.M.: Overview of the TREC 2011 medical records track. In: TREC 2011, Gaithersburg, Maryland, NIST (2011)
23. Voorhees, E.M., Hersh, W.: Overview of the TREC 2012 medical records track. In: TREC 2012, Gaithersburg, Maryland, NIST (2012)