

Integrating Social Features and Query Type Recognition in the Suggestion Track of CLEF 2015 Social Book Search Lab

Shih-Hung Wu^{1*}, Yi-Hsiang Hsieh¹, Liang-Pu Chen², Tsun Ku²

¹Chaoyang University of Technology, Taiwan, R.O.C
{shwu(*Contact author), s10027006,}@cyut.edu.tw

²Institute for Information Industry, Taiwan, R.O.C
eit@iii.org.tw, cujing@gmail

Abstract. The Social Book Search (SBS) Lab is part of CLEF 2015 lab series. This is the third time that the CYUT CSIE team attends the SBS track. Based on a full-text search engine, we build a social feature re-ranking system and introduce more knowledge on understanding the queries. We defined a set of rules to filtering out unnecessary books from the recommendation list. The official run results show that the system performance is improved from our previous system.

Keywords: Query type recognition, social features, social book search

1 Introduction

The paper reports our system in the suggestion track of CLEF 2015 Social Book Suggestion (SBS) [10]. This is the third time that we attend the SBS track since 2013 INEX [7]. Based on our social feature re-ranking system [1], we improve our system by involving some knowledge on understanding the queries.

We believe that the result of traditional information retrieval technology is not enough for the users who need more personal recommendation in the SBS task. Recommendation from other users are more appealing; it might contain more personal feelings and cover more subtle reasons that traditional information retrieval system cannot cover. Our system integrates the social feature into the traditional information retrieval technology to give better recommendation on books. In this task, user-generated metadata is used as the social feature.

According to our observation on the topics in the previous INEX SBS Track, we found that queries can be separated into different types. Simply treating the keywords in the topic as search terms will not get good results. Some queries require higher level of knowledge to deal with. System needs to understand the information need behind the keyword, for example, the knowledge on the types of literature. We analysis the topics and find several types in them. Due to the time limitation, we only implement a module to recognize one special type of topics and a filtering module to modify the recommendation result.

The structure of this paper is as follows. Section 2 is the data set description, sec-

tion 3 shows our architecture and the details of our method, section 4 is the experiment results, and final section gives conclusions and future works.

2 Dataset

2.1 Collection

The document collection in this task is provided by the CLEF 2015 Social Book Suggestion track. The documents are the XML format metadata of about 2.8 million books and the data size is 25.9GB. These documents are collected from Amazon.com and LibraryThing [2]. The XML tags used in the data set is listed in Table 1.

Table 1.All the XML tag [2]

tag name			
book	similarproducts	title	imagecategory
dimensions	Tags	edition	name
reviews	Isbn	dewey	role
editorialreviews	Ean	creator	blurber
images	Binding	review	dedication
creators	Label	rating	epigraph
blurbers	Listprice	authorid	firstwordsitem
dedications	manufacturer	totalvotes	lastwordsitem
epigraphs	numberofpages	helpfulvotes	quotation
firstwords	publisher	date	seriesitem
lastwords	Height	summary	award
quotations	Width	editorialreview	browseNode
series	Length	content	character
awards	Weight	source	place
browseNodes	readinglevel	image	subject
characters	releasedate	imageCategories	similarproduct
places	publicationdate	url	tag
subjects	Studio	data	

2.2 Test Topic

Topics provided by CLEF 2015 Social Book Suggestion track are collected from LibraryThing. A topic describes the information needed for a user. Figure 1 and Figure 2 give partial view of an example, the XML tags used are : <topic id>, <title>, <mediated_query>, <group>, <narrative>, <catalog>, <book>, <LT_id>, <entry_date>, and <rating>. Where title means the title of a post on LibraryThing forum and narrative is the content of the post. While mediated_query is added as an interpretation of the query. Group means the user group in the forum of the user who post this query.

```
<topics>
```

```

<topic id="1196">
  <title>The Best Peace Corps Novel</title>
  <mediated_query>books about work for Peace Corps </mediate
d_query>
  <group>Returned Peace Corps Volunteer Readers</group>
  <narrative>
    I'm looking for people's concept of what is
    the best novel for the Peace Corps Volunteer - pre, during, o
    r post service. This could be a novel that typifies life in th
    e country of service. It could be a novel that typifies the wo
    rk volunteers do. It could be a novel that makes for the perfe
    ct reading while in service. Anything will do, just give rea
    sons. It might lead other PCVs/RPCVs to interesting reading.
    Let's try novels, and then head into non-fiction later... I'll
    start: I could not have survived my 2 years of service if I
    had not read Chingiz Aitmitov's The Day Lasts More than A Hun
    dred Years and Bulgakov's The Master and Margarita . They re
    ally made most of my concerns about my own sanity living in th
    e crumbling remnants of Soviet Central Asia vanish into vapor,
    as I was able to learn that not only was surreality the norm
    for this part of the world but also my own preconceptions abou
    t the concrete, rational world that I thought I knew might be
    questionable. </narrative>

```

Figure 1. A topic example in CLEF 2015 Social Book Suggestion track

```

<examples>
  <example>
    <LT_id>120241</LT_id>
    <hasRead>yes</hasRead>
    <sentiment>positive</sentiment>
  </example>
  <example>
    <LT_id>10151</LT_id>
    <hasRead>yes</hasRead>
    <sentiment>positive</sentiment>
  </example>
</examples>
<catalog>
  <book>
    <LT_id>42437</LT_id>
    <entry_date>2006-07</entry_date>
    <rating>10.0</rating>
    <tags></tags>
  </book>
</book>

```

```

    <LT_id>1270696</LT_id>
    <entry_date>2006-07</entry_date>
    <rating>10.0</rating>
    <tags></tags>
  </book>
  ...
</catalog>
</topic>

```

Figure 2. A topic example in CLEF 2015 Social Book Suggestion track (Continued)

3 CYUT CSIE System Methodology

3.1 System Architecture

Figure 3 shows our system architecture. The pre-processing modules include stop words filtering, and stemming, both modules are provided by Lucene [6]. After the preprocessing, our system builds index for retrieval. The results of content-based retrieval will be re-ranked as the final results according to the social features.

3.2 Indexing and Query

The index and search engine in use is the Lucene system, which is an open source full text search engine provided by Apache software foundation. Lucene is written in JAVA and can be called easily by JAVA program to build various applications.

Table 1 shows all the tags of the book metadata. According to Bogers and Larsen [3], there are 19 tags more useful in the social book search. They are <isbn>, <title>, <publisher>, <editorial>, <creator>, <series>, <award>, <character>, <place>, <blurber>, <epigraph>, <firstwords>, <lastwords>, <quotation>, <dewey>, <subject>, <browseNode>, <review>, and <tag>. Our system also focuses on the same 19 tags.

In the pre-processing step, the content in the <dewey> tag is restored to strings according to the 2003 list of Dewey category descriptions [9] to make string matching easier. For example: <dewey>004</dewey> will be restored to <dewey>Data processing Computer science</dewey>. The content of <tag> is also expanded according to the count number to emphasize its importance. For example: <tag count="3">fantasy</tag> will be expanded as <tag>fantasy fantasy fantasy</tag>. In addition to the 19 tags, our system also indexes the content of <review> as independent indexes files and names it as reviews.

Fig.1 and 2 shows all the XML tags of the query topics. According to Koolen et al. [4], an Indri [5] based system using all the contents of <Title>, <Query>, <Group>, and <Narrative> as query terms will give better result. We also use the contents of the four tags as our system input queries.

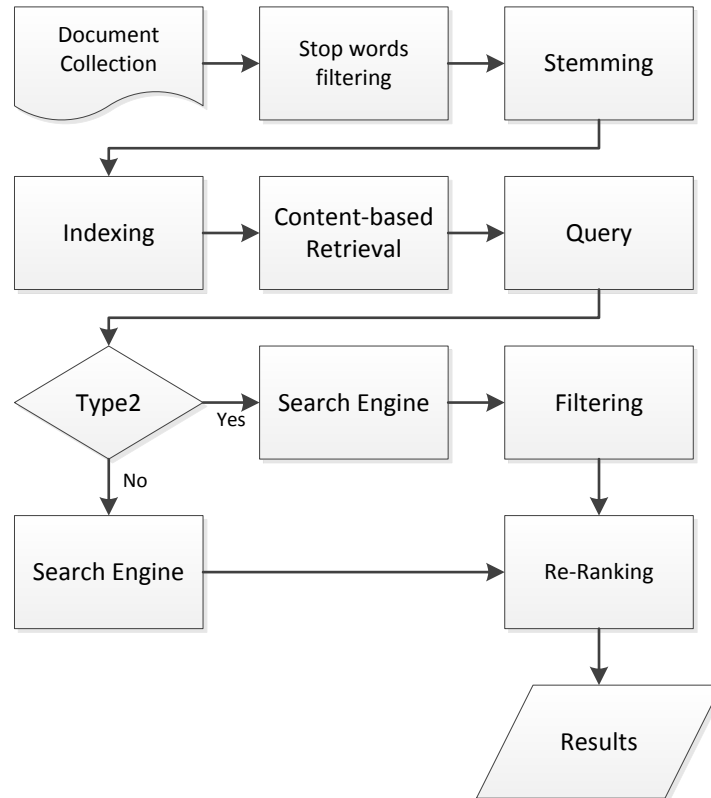


Figure 3. System architecture

```

<topic id="76778">
  <title>Russian Serfdom Suggestions</title>
  <mediated_query>Russian serfdom </mediated_query>
  <group>History Readers: Clio's (Pleasure?) Palace</group>
  <narrative>I'm reading Flashman At The Charge right now
    and Russian serfdom is a prominent feature. Any
    one have any good suggestions to learn more abo
    ut this aspect of Russian history during the Ts
    ars? I'm looking for a Gulag: A History about
    serfdom. Thanks! </narrative>
  <examples>

```

Figure 4. A type2 query example that we defined in 2015 SBS track

3.3 Type2 Query Recognition and Result Filtering

According to our observation on the topics in INEX 2012 SBS Track, we find that there are some queries that are different from others, we call them the Type2 queries [11]. Type2 queries are the queries that contain the names of some books that the original users want to find similar ones. Therefore, the books in the topics should not be part of the recommendation. Since the book names are given explicitly, our system originally will find exactly the same books as the top recommendation. To recognize type2 queries, we define a list of phrases to identify such queries and filter out the books in the queries from the recommendation lists. The phrases are listed in the appendix in the previous paper [11]. Figure 4 gives an example of Type2 queries taken from INEX 2013 SBS topics, in which contains a key phrase “I’m reading”. We find that there are 174 queries in the INEX 2013 SBS track that can be classified as Type2 queries. Therefore, we add a module in our system to identify the Type2 queries and filtering out the books mentioned in the topics.

3.4 Re-ranking

The Re-ranking part is similar to that in our previous work [1]. We integrate the user-generated metadata into the traditional content-based search result by re-ranking the results. The social features are provided by the amazon users, and our system use them to give more weight on certain books. Three numbers are available:

- User rating: users might evaluate a book from 1 to 5, the higher the better.
- Helpful vote: other users might endorse one comment by voting it as helpful.
- Total vote: the total number of helpful or not.

We designed 3 different ways to use these social features in re-ranking.

1) User rating method

Increase the weight of content-based retrieval result by adding the summation of user rating. As shown in formula (1):

$$\text{Score}_{\text{re-ranked}}(i) = \alpha * \text{Score}_{\text{org}}(i) + (1 - \alpha) * \text{Score}_{\text{user rating}}(i) \quad (1)$$

2) Average User rating method

Increase the weight of content-based retrieval result by adding the average of user rating. As shown in formula (2):

$$\text{Score}_{\text{re-ranked}}(i) = \text{Score}_{\text{org}}(i) + \text{Score}_{\text{average user rating}}(i) \quad (2)$$

3) Weights User rating method

Increase the weight of content-based retrieval result by adding the book which gets more helpful votes. As shown in formula (3) and (4):

$$\text{Score}_{\text{Weights User Rating}} = \text{User rating} * \frac{\text{helpfulvote}}{\text{totalvote}} \quad (3)$$

$$\text{Score}_{\text{re-ranked}}(i) = \alpha * \text{Score}_{\text{org}}(i) + (1 - \alpha) * \text{Score}_{\text{Weights User Rating}}(i) \quad (4)$$

3.5 Find the Best α Value by Experiment

Since there is no theoretical reference on how to set the α value, in our official runs, the value is selected via a series experiments that we conduct on the 2013 dataset. Table 2 shows the results, we find that the system gets the best result when α is 0.95.

Table 2. Experimental Result for different α on 2013 data set

A	P@10	MAP
0.50	0.0221	0.0193
0.60	0.0221	0.0193
0.70	0.0224	0.0195
0.80	0.0226	0.0196
0.90	0.0237	0.0204
0.95	0.0245	0.0220

4 Experimental results

In the official evaluation, we sent four runs. We use four fields in the topics as query terms, and we filter out some book candidates for all the type2 queries. The configuration of each run is as follows.

- Run 1, the CSIE - 0.95AverageType2QTGN, re-ranking with Average User Rating.
 - Run 2, the CSIE - Type2QTGN: without re-ranking.
 - Run 3, the CSIE - 0.95RatingType2QTGN, re-ranking with User Rating.
 - Run 4, CSIE - 0.95WRType2QTGN, Re-ranking with Weights User Rating.
- According to Table 2, the parameter α is 0.95 for best result in the runs with re-ranking.

Table 3 shows the official evaluation results of our four runs. Among them the CSIE - 0.95AverageType2QTGN run gives the best NDCG@10 [8] result, while the CSIE - Type2QTGN run gives similar result on NDCG@10 but give better result on MAP and R@1000. The other two runs give poorer results might due to technical errors. Comparing to the 2013 INEX SBS results in Table 5, our system performance improved significantly. However, comparing to the result of INEX SBS 2014 in Table 4, our system performance decreased.

Table 3. Official evaluation results in 2015 SBS

Run	nDCG@10	MRR	MAP	R@1000	Profiles
CSIE - 0.95AverageType2QTGN	0.082	0.194	0.050	0.319	no
CSIE - Type2QTGN	0.080	0.191	0.052	0.325	no
CSIE - 0.95RatingType2QTGN	0.032	0.113	0.019	0.214	no
CSIE - 0.95WRType2QTGN	0.023	0.072	0.015	0.216	no

Table 4. Official evaluation results in 2014 INEX SBS

Run	nDCG@10	MRR	MAP	R@1000
CYUT - Type2QTGN	0.119	0.246	0.086	0.340
CYUT -	0.119	0.243	0.085	0.332

0.95AverageType2QTGN				
CYUT - 0.95RatingType2QTGN	0.034	0.101	0.021	0.200
CYUT - 0.95WRType2QTGN	0.028	0.084	0.018	0.213

Table 5. Official evaluation results in 2013 INEX SBS

<i>Run</i>	<i>nDCG@10</i>	<i>P@10</i>	<i>MRR</i>	<i>MAP</i>
Run1.query.content-base	0.0265	0.0147	0.0418	0.0153
Run2.query.Rating	0.0376	0.0284	0.0792	0.0178
Run3.query.RA	0.0170	0.0087	0.0352	0.0107
Run4.query.RW	0.0392	0.0287	0.0796	0.0201
Run5.query.reviwes.content-base	0.0254	0.0153	0.0359	0.0137
Run6.query.reviews.RW	0.0378	0.0284	0.0772	0.0165

5 Conclusions and Future work

This paper reports our system and result in CLEF 2015 Social Book Suggestion track. We sent four runs and the formal run results are list in Table 3. In the four runs, the CSIE - 0.95AverageType2QTGN run gives best nDCG@10, which is searching with content-based search engine, applying a set of filtering rules based on a list of key phrase and re-ranking with Average User Rating. In the future, we will implement more modules with literature knowledge on the writers, genre of books, geometric categories of the publishers, and temporal categories of the authors that can deal with the special cases in the topics.

From this year, user profiles are available, which can be used to give better recommendation. A system might use the user profiles to expand the queries or to suggest more books that the user read before for other similar users. Outside resources might also be used to expand the queries. For example, a system might check Wikipedia to find more authors of the books in the same genre, and make better recommendation. Books that won some awards might also be a good list for recommendation.

Acknowledgement

“This study is conducted under the "Online and Offline integrated Smart Commerce Platform(2/4)" of the Institute for Information Industry which is subsidized by the Ministry of Economy Affairs of the Republic of China .

References

1. Wei-Lun Xiao, Shih-Hung Wu, Liang-Pu Chen, Hung-Sheng Chiu, and Ren-Dar Yang, “Social Feature Re-ranking in INEX 2013 Social Book Search Track”, CLEF 2013 Evaluation Labs and Workshop Online Working Notes, 23 - 26 September, Valencia, Spain.
2. Marijn Koolen, Gabriella Kazai, Jaap Kamps, Michael Preminger, Antoine Doucet, and Monica Landoni, “Overview of the INEX 2012 Social Book Search Track”, INEX’12

- Workshop Pre-proceedings,P.77-P.96,2012.
3. Toine Bogers and Birger Larsen, “RSLIS at INEX 2012: Social Book Search Track”, INEX'12 Workshop Pre-proceedings,P.97-P.108,2012.
 4. Marijn Koolen, Hugo Huurdeman and Jaap Kamps, “Comparing Topic Representations for Social Book Search”, CLEF 2013 Evaluation Labs and Workshop Online Working Notes, 23 - 26 September, Valencia – Spain.
 5. T. Strohmman, D. Metzler, H. Turtle, and W. B. Croft, “Indri: a language-model based search engine for complex queries”, In Proceedings of the International Conference on Intelligent Analysis, 2005.
 6. Lucene, <https://lucene.apache.org>
 7. Marijn Koolen, Gabriella Kazai, Michael Preminger, and Antoine Doucet, “Overview of the INEX 2013 Social Book Search Track”, CLEF 2013 Evaluation Labs and Workshop Online Working Notes, 23 - 26 September, Valencia – Spain.
 8. Järvelin, K., Kekäläinen, “J.: Cumulated Gain-based Evaluation of IR Techniques”, ACM Transactions on Information Systems 20(4) (2002) 422–446.
 9. 2003 list of Dewey category descriptions, <https://www.library.illinois.edu/ugl/about/dewey.html>
 10. CLEF 2015 Social Book Search Track, <http://social-book-search.humanities.uva.nl/#/suggestion>
 11. Shih-Hung Wu, Pei-Kai Liao, Hua-Wei Lin, Li-Jen Hsu, Wei-Lun Xiao, Liang-Pu Chen, Tsun Ku, and Gwo-Dong Chen Query Type Recognition and Result Filtering in INEX 2014 Social Book Search Track
 12. Marijn Koolen, Toine Bogers, Gabriella Kazai, Jaap Kamps, and Michael Preminger Overview of the INEX 2014 Social Book Search Track