

An Approach to Structured Retrieval Based on the Extended Vector Model

Carolyn J. Crouch
Department of Computer Science
University of Minnesota Duluth
Duluth, MN 55812
(218) 726-7607
ccrouch@d.umn.edu

Sameer Apte
Department of Computer Science
University of Minnesota Duluth
Duluth, MN 55812
(218) 726-7607
apte0002@d.umn.edu

Harsh Bapat
Persistent Systems Pvt. Ltd.
Pune, India 411016
+91 (20) 567 8900
harsh_bapat@persistent.co.in

ABSTRACT

In this paper, we describe our approach to XML retrieval, which is based on the extended vector space model initially proposed by Fox [5]. The current implementation of our system and results to date are reported. The basic functions are performed using the Smart experimental retrieval system. Early results confirm the viability of the extended vector space model in this environment.

1. INTRODUCTION

When we began our work with INEX last year, our goal was to confirm the utility of Salton's vector space model [10] in its extended form for XML retrieval. Long familiarity with Smart [9] and its capabilities led us to believe that it could be used for this purpose. Our approach was described in the proceedings of last year's workshop [3]. Much initial effort was spent on the translation of documents and topics from XML to internal Smart format and the subsequent translation of results back into INEX format. When we reported our results in [3], our system was still in a very rudimentary stage.

In 2002, we had an idea and began implementation. During the past year, we have built upon and extended that work. We now have an operational system. For the sake of clarity, a brief overview follows.

1.1 Background

Everyone involved in information retrieval is familiar with the vector space model, wherein documents and queries are represented as weighted term vectors. The weight assigned to a term is indicative of the contribution of that term to the meaning of the document. Very commonly, *tf-idf* weights [11] or some variation thereof [12] are used. The similarity between vectors (e.g., document and query) is represented by the mathematical similarity of their corresponding term vectors.

In 1983, Fox [5] proposed an extension of the vector space model—the so-called *extended vector space model*—to allow for the incorporation of objective identifiers with content identifiers in the representation of a document. An extended vector can include different classes of information about a document, such as author name, publication date, etc., along with content terms. In this model, a document vector consists of a set of subvectors, where each subvector represents a different class of

information (i.e., concept class or c-type). Our current representation of an XML document/topic consists of 18 c-types (i.e., *abs*, *ack*, *articl_au_fnm*, *article_au_snm*, *atl*, *au_aff*, *bibl_atl*, *bibl_au_fnm*, *bibl_au_snm*, *bibl_ti*, *ed_aff*, *ed_intro*, *kwd*, *rname*, *st*, *ti*, *pub_yr*, *bdy*) as defined in INEX guidelines. Subjective subvectors are those with a body of text associated with them (i.e., *abs*, *ack*, *atl*, *bibl_atl*, *bibl_ti*, *ed_intro*, *kwd*, *bdy*). Similarity between extended vectors is calculated as a linear combination of the similarities of the corresponding subvectors.

Use of the extended vector model for document retrieval normally raises at least two problems: the construction of the extended search request [4, 6] and the selection of the coefficients for combining subvector similarities. For XML retrieval, the CO query in particular can be roughly translated into extended vector form by distributing the keywords across the subjective subvectors. (CAS queries are more difficult; we are working on automating this process.) The second problem—the weighting of the subvectors themselves—remains open to investigation. Another issue of some interest here is the weighting of terms within the subvectors. (We have produced some useful results in relation to the term weighting issue; our work on the weighting of subvectors is promising but not well developed. In any case, subvector weighting is unlikely to have a measurable effect within the large INEX window.)

The extended vector capability of Smart appeared to us well suited for XML with respect to the retrieval of documents. But there is no facility for retrieving at the element level (or at various levels of granularity), which is a requirement of INEX tasks. We are interested in determining the feasibility of incorporating the functionality (i.e., flexibility and granularity) required for XML retrieval within the extended vector environment. We are currently investigating methods that have been suggested by others (e.g., Grabs and Schek [7, 8]). However, more work is necessary before conclusions can be drawn.

1.2 System Description

Our system handles the processing of XML text as follows:

- (1) The documents are parsed using a simple XML parser available on the web. Each of our 18 c-types is now identifiable in terms of its XML path.

- (2) The documents and queries are translated into Smart format and indexed by Smart as extended vectors. (The results reported in this paper are all based on an indexing which considers the body of the document as a single entity; i.e., paragraphs and sections, for example, are not recognized.)
- (3) Retrieval takes place by running the queries against the indexed collection. The result is a list of articles ordered by decreasing similarity to the query. (A number of term weighting schemes are available through Smart.)
- (4) For each query, the top 100 articles are converted back into INEX format and reported.

The retrieval itself is straight-forward. The only variation is the splitting of certain CAS queries into separate portions which are then run in parallel to ensure that the elements retrieved meet the specified criteria. See Section 2.2 for an example of this type.

2. EXPERIMENTS

In the following sections, we describe the experiments performed with respect to the processing of the CO and CAS topics, respectively. In all cases, we used only the topic title and keywords as search words in query construction. As indicated previously, this year's effort focused on producing a working system—by our definition, a system that returns competitive results with respect to at least some INEX task(s). To demonstrate that our system is functional, we first processed the INEX 2002 topics (under the original *inex_eval*) to compare our results to those already reported. We then processed the 2003 topics. The results are all reported here.

2.1 Using CO Topics

Our first task is to formulate the CO topic in extended vector form. Of the 18 c-types composing the extended vector, 8 contain subjective identifiers (i.e., *abs*, *ack*, *atl*, *bibl_atl*, *bibl_ti*, *ed_intro*, *kwd*, *bdy*). The extended vector topic is formed by associating the search words of the topic with each of these 8 c-types. The remaining c-types contain objective identifiers and are not used in formulating CO queries. Our more interesting experiments are discussed briefly below. (See [1] for details.) The subvectors are equally weighted in all these cases.

2.1.1 2002 Topics

Our term weighting experiments include:

Tuned *Lnu-ltu* Term Weighting: In this experiment, we tuned the collection as indicated by Singhal, *et. al.*, in [13]. Results under generalized quantization were 0.065 whereas strict quantization produced 0.095.

Augmented *tf-idf* (*atc*) Term Weighting: 2002 topics under generalized quantization produced an average precision of 0.033.

Retrieval at the Element Level: In this experiment, we used indexings of the collection at the paragraph and section levels in addition to the article level. (Untuned or estimated *Lnu-ltu* weights were used in these early experiments.) For each query, the rank-ordered lists were sorted and the top 100 elements reported. Average precision was 0.042 under generalized quantization.

2.1.2 2003 Topics

Our 2003 CO submission was based on parameters that produced the best results for 2002 CO topics, i.e., *Lnu-ltu* term weighting with equal subvector weights. The recall-precision graphs for 2003 CO topics under the revised *inex_eval* are given below in Figures 1 and 2. The results under *inex_eval_ng* (overlap ignored) are shown in Figures 3 and 4. Corresponding results under *inex_eval_ng* (overlap considered) are shown in Figures 5 and 6.

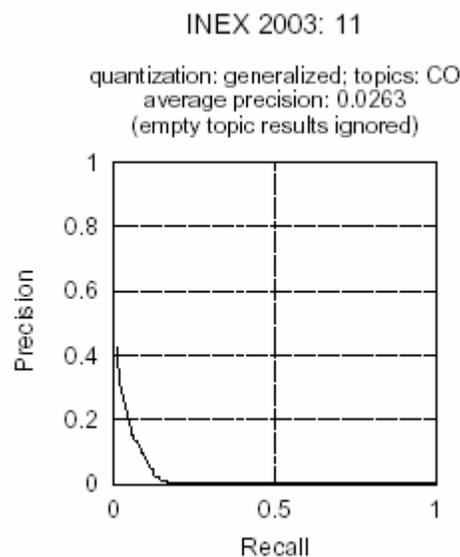


Figure 1. Recall-precision for CO, Gen

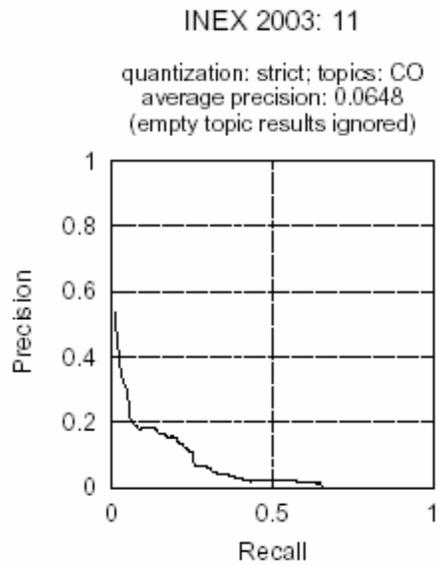


Figure 2: Recall-precision for CO, Strict

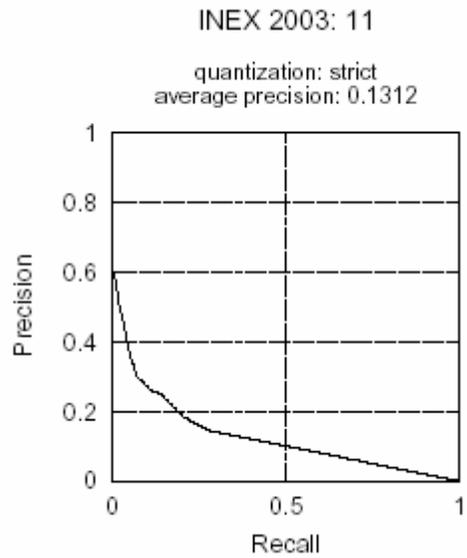


Figure 4: Recall-precision for CO, Strict under ng (overlap ignored)

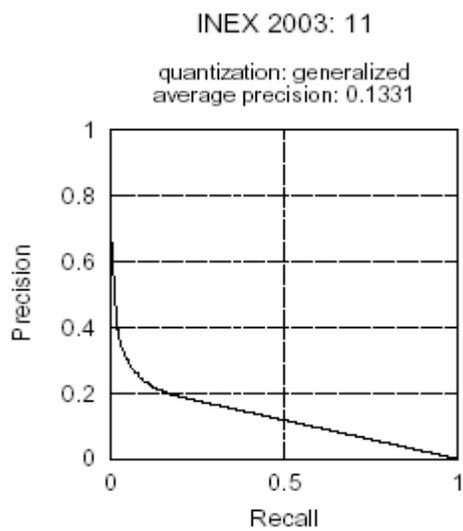


Figure 3: Recall-precision for CO, Gen under ng (overlap ignored)

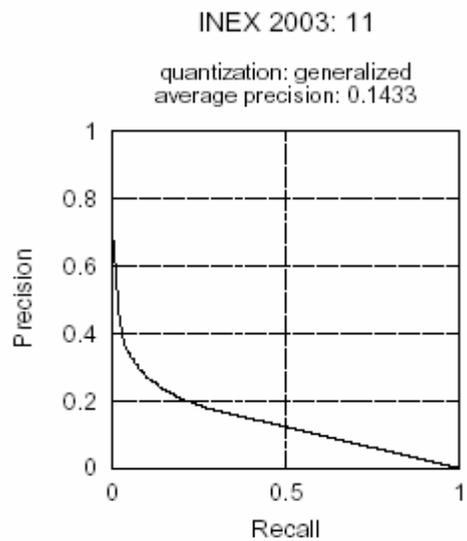


Figure 5: Recall-precision for CO, Gen under ng (overlap considered)

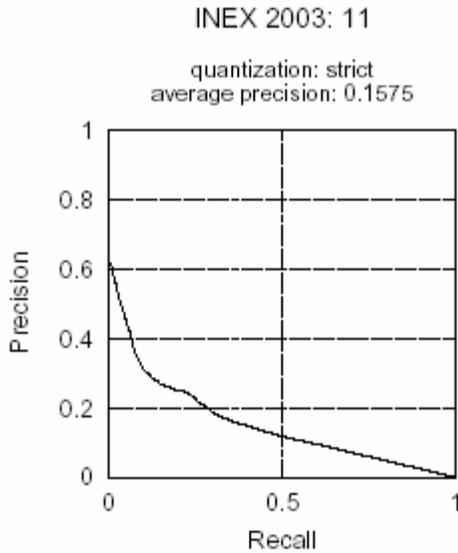


Figure 6: Recall-precision for CO, Strict under *ng* (overlap considered)

To recap: Our results for 2003 INEX CO topics are on the whole good, ranking in the top 10 of the 4 evaluations (Figures 3, 4, 5, and 6) under *inex_eval_ng*. Yet although we were able to produce decent results for the 2002 CO topics under the original *inex_eval*, our results for the 2003 CO topics under the revised *inex_eval* fall far from the top. We are still assessing the causes.

2.2 Using CAS Topics

We were able to formalize the extended vector CO topics fairly easily. The extended vector CAS topic formulations, on the other hand, present more of a challenge. Direct use of the extended vector model does not guarantee that each keyword will occur in the specified context. To effect this result, we currently split certain CAS queries into separate portions which are then run in parallel to ensure that the elements retrieved meet the specified criteria. Consider, for example, the title section of CAS query 8:

```
<title>
  <te>article</te>
  <cw>ibm</cw><ce>fm/aff</ce>
  <cw>certificates</cw><ce>bdy/sec</ce>
</title>
```

In this case, the query is to return a ranked list of articles as specified by the target element `<te>`. The narrative specifies that the body or sections of relevant documents should contain information about the use of certificates for authenticating users on the Internet. And since the context of the content word *ibm* is *fm/aff*, the author(s) of those documents must be affiliated with IBM. Thus the query

should retrieve only those articles on the use of certificates whose author(s) are affiliated with IBM. To guarantee that the system returns *only* those articles, we split the query into two parallel queries as follows:

Q1: `<cw>ibm</cw><ce>fm/aff</ce>`

Q2: `<cw>certificates</cw><ce>bdy/sec</ce>`

Affiliation and section are two different c-types. So query 1 searches for documents containing the objective identifier *ibm* in the affiliation subvector. Query 2 seeks articles whose body or section(s) contain the term *certificate*. Smart returns a ranked list of documents for both queries. The intersection of these lists is the final, ranked list of documents returned. This approach—the splitting of a query into parts—is a first step in the process of using objective ctypes to filter results appropriately.

This year we experimented with different term weighting schemes for CAS topics. We performed these experiments first on the 2002 topics. Equal subvector weighting was applied in each case. Experiments performed during the past year using the INEX 2002 queries are described briefly below. (See [2] for details.) Evaluation for these topics was performed through the original *inex_eval*.

2.2.1 2002 Topics

Untuned *Lnu_ltu* Term Weighting: All subvectors are weighted in this fashion. Average precision was 0.179 under generalized and 0.222 under strict quantization.

Lnu_ltu (for subjective subvectors) and *nnn* (for objective subvectors) Term Weighting: Here we used simple term frequency weights (*nnn*) for the objective subvectors combined with *Lnu_ltu* weights for the subjective subvectors. Average precision was 0.187 under generalized and 0.235 under strict quantization.

Augmented tf-idf (*atc*) Term Weighting: All subvectors were weighted with *atc* weights. Average precision was 0.194 and 0.238 under generalized and strict quantization, respectively.

Augmented tf-idf (*atc*—for subjective subvectors) and *nnn* (for objective subvectors) Term Weighting: These weights returned an average precision of 0.192 under generalized and 0.243 under strict quantization.

All of these results rank in the top 10 when compared to the best case results reported for INEX 2002 topics.

2.2.2 2003 Topics

Our 2003 submission used *atc* term weighting for all subvectors with equal subvector weights. Due to the exigencies of the academic schedule, we were able to submit only under VCAS. Results await availability of the corresponding INEX evaluation package, but we do not expect them to be useful at this point. We need to modify our methods so that the appropriate filters are applied

before results are returned.

During the past year, we produced a working system. An overview of our results may be seen in Table 1. The column labeled UMD (for University of Minnesota Duluth) presents our results, which may be compared with the best result reported for that task (in the INEX column).

3. CONCLUSIONS

In 2003, our efforts were directed at producing a working system for structured retrieval based on the extended vector model. In our view, this year's results have demonstrated the viability of such an approach. However, structured retrieval requires additional capabilities beyond the scope of normal vector-based systems, and thus the question remains. Is our model—the extended vector model—able to support the functionality required in this environment?

Our system is still in an early stage of development. The issue of term weighting has now become clearer; the weighting of the subvectors themselves is still an open question. The major challenge is to develop a method of returning results at the element level, i.e., to retrieve at the desired level of granularity. Our plans include further investigation of the methods of others along with the development of an approach that may be better suited to our own environment. Another major focus is the development of appropriate techniques for handling CAS topics effectively.

Table 1. Comparison of Best Case Avg Precision for CO Topics

	UMD		INEX	
	gen	strict	gen	strict
'02 Topics	0.0650	0.0950	0.0700	0.0880
'03 Topics: inex_eval	0.0263	0.0648	0.1036	0.1214
'03 Topics: inex_eval_ng*	0.1331	0.1312	0.1783	0.1857
'03 Topics: inex_eval_ng**	0.1433	0.1575	0.1542	0.1584

* overlap ignored; ** overlap considered

4. REFERENCES

[1] Apte, S. Adapting the extended vector space model for content-oriented XML retrieval. Master's Thesis, Dept. of Computer Science, University of Minnesota Duluth (2003). www.d.umn.edu/~ccrouch/publications.html

[2] Bapat, H. Adapting the extended vector space model for structured XML retrieval. Master's Thesis, Dept. of Computer Science, University of Minnesota Duluth (2003). www.d.umn.edu/~ccrouch/publications.html

[3] Crouch, C., Apte, S., and Bapat, H. Using the extended vector model for XML retrieval. In Proc of the First Workshop of the Initiative for the Evaluation of XML Retrieval (INEX), (Schloss Dagstuhl, 2002), 99-104.

[4] Crouch, C., Crouch, D. and Nareddy, K. The automatic generation of extended queries. In Proc. of the 13th Annual International ACM SIGIR Conference, (Brussels, 1990), 369-383.

[5] Fox, E. A. Extending the Boolean and vector space models of information retrieval with p-norm queries and multiple concept types. Ph.D. Dissertation, Department of Computer Science, Cornell University (1983).

[6] Fox, E., Nunn, G. and Lee, W. Coefficients for combining concept classes in a collection. In Proc. of the 11th Annual International ACM SIGIR Conference, (Grenoble, 1988), 291-307.

[7] Grabs, T. and Schek, H. Generating vector spaces on-the-fly for flexible XML retrieval. In Proc of the ACM SIGIR Workshop on XML and Information Retrieval, (Tampere, Finland, 2002), 4-13.

[8] Grabs, T. and Schek, H. ETH Zurich at INEX: Flexible information retrieval from XML with PowerDB-XML. INEX 2002 Workshop Proceedings, (Dortland, 2002), 35-40.

[9] Salton, G. Automatic information organization and retrieval. Addison-Wesley, Reading PA (1968).

[10] Salton, G., Wong, A., and Yang, C. S. A vector space model for automatic indexing. Comm. ACM 18, 11 (1975), 613-620.

[11] Salton, G. and Buckley, C. Term weighting approaches in automatic text retrieval. In IP&M 24, 5 (1988), 513-523.

[12] Singhal, A. AT&T at TREC-6. In The Sixth Text RETrieval Conf (TREC-6), NIST SP 500-240 (1998), 215-225.

[13] Singhal, A., Buckley, C., and Mitra, M. Pivoted document length normalization. In Proc. Of the 19th Annual International ACM SIGIR Conference, (Zurich, 1996), 21-19.