

Preferential Infinitesimals for Information Retrieval

Maria Chowdhury, Alex Thomo, and William W. Wadge

Department of Computer Science, University of Victoria, Canada
{mwchow, thomo, wwadge}@cs.uvic.ca

Abstract. In this paper, we propose a preference framework for information retrieval in which the user and the system administrator are enabled to express preference annotations on search keywords and document elements, respectively. Our framework is flexible and allows expressing preferences such as “ A is infinitely more preferred than B ,” which we capture by using *hyperreal numbers*. Due to the widespread of XML as a standard for representing documents, we consider XML documents in this paper and propose a consistent preferential weighting scheme for nested document elements. We show how to naturally incorporate preferences on search keywords and document elements into an IR ranking process using the well-known TF-IDF ranking measure.

1 Introduction

In this paper, we propose a framework for preferential information retrieval by incorporating in the document ranking process preferences given by the user or the system administrator. Namely, in our proposed framework, the user has the option of weighting the search keywords, whereas the system administrator has the option of weighting structural elements of the documents. We address both facets of preferential weighting by using hyperreal numbers, which form a superset of the real numbers, and in our context, serve the purpose of specifying natural preferences of the form “ A is infinitely more preferred than B .”

Keyword Preferences. To illustrate preferences on keywords, suppose that a user wants to retrieve documents on research and techniques for “music-information-retrieval.” Also, suppose that the user is a fan of Google technology. As such, this user would probably give to a search engine the keywords:

music-information-retrieval, google-search, google-ranking.

It is interesting to observe that if the user specifies these keywords in Google, then she gets a list of only *three*, low quality, pages. What happens is that the true, highly informative pages about “music-information-retrieval” are lost (or insignificantly ranked) in the quest of trying to serve the “google-search” and “google-ranking” keywords. Unfortunately, in Google and other search engines, the user cannot explicitly specify her real preferences among the specified keywords. In this example, what the user needs is a mechanism for saying that “music-information-retrieval” is of primary importance or *infinitely* more important than “google-search” and “google-ranking,” and thus, an informative page about “music-information-retrieval” should be retrieved and highly ranked even if it does not relate to Google technologies.

Structural Preferences. The other facet of using preferential weights is for system administrators to annotate structural parts of the documents in a given corpus. In practice, most of the documents are structured, and often, certain parts of them are more important than others. While our proposed ideas can be applied on any corpus of structured documents, due to the wide spread of XML as a standard for representing documents, we consider in this paper XML documents which conform to a given schema (DTD). In the same spirit as for keyword preferences, we will use hyperreal weights to denote the importance of different elements in the schema and documents.

To illustrate preferences on structural parts of documents, suppose that we have a corpus of documents representing research papers, and a user is searching for a specific keyword. Now, suppose that the keyword occurs in the *title* element of one paper and in the *references* element of another paper. Intuitively, the paper having the keyword in the *title* should be ranked higher than the paper containing the keyword in the *references* element as the title of a paper usually bears more representative and concise information about the

paper than the reference entries do. In fact, one could say that terms in the title (and abstract) are *infinitely* more important than terms in the references entries as the latter might be there completely incidental.

While weighting of certain parts of documents has been considered and advocated in the folklore (cf. [6, 9]), to the best of our knowledge there is no work dealing with inferring a consistent weighting scheme for nested XML elements based on the weights that a system administrator gives to DTD elements. As we explain in Section 4, there are tradeoffs to be considered and we present a solution that properly normalizes the element weights producing values which are consistent among sibling elements and never greater than the normalized weight of the parent element, thus respecting the XML hierarchy.

Contributions. Specifically, our contributions in this paper are as follows.

1. We propose using hyperreal numbers (see [7, 8]) to capture both “quantitative” and “qualitative” user preferences on search keywords. The set of hyperreal numbers includes the real numbers which can be used for expressing “quantitative” preferences such as, say “*A* is twice more preferred than *B*,” as well as *infinitesimal* numbers, which can be used to express “qualitative” preferences such as, say “*A* is infinitely more preferred than *B*.” We argue that without such qualitative preferences there is no guarantee that an IR system would not override user preferences in favor of other measures that the system might use.
2. We extend the ideas of using hyperreal numbers to annotating XML (DTD) schemas. This allows system administrators to preferentially weight structural elements in XML documents of a given corpus. We present a normalization method which produces consistent preferential weights for the elements of any XML document that complies to an annotated DTD schema.
3. We adapt the well-known TF-IDF ranking in IR systems to take into consideration the preferential weights that the search keywords and XML elements can have. Our extensions are based on symbolic computations which can be effectively computed on expressions containing hyperreal numbers.
4. We present (in the appendix) illustrative practical examples which demonstrate the usefulness of our proposed preference framework. Namely, we use a full collection of speeches from the Shakespeare plays, and a diverse XML collection from INEX ([15]). In both these collections, we observed a clear advantage of our preferential ranking over the ranking produced by the classical TF-IDF method. We believe that these results encourage incorporating both quantitative and (especially) qualitative preferences into other ranking methods as well.

Organization. The rest of the paper is organized as follows. In Section 2, we give an overview of hyperreal numbers and their properties. In Section 3, we present hyperreal preferences for annotating search keywords. In Section 4, we propose annotated DTDs for XML documents and address two problems for consistent weighting of document elements. In Section 5, we show how to extend the TF-IDF ranking scheme to take into consideration the hyperreal weights present in the search keywords and document elements. In Appendix, we present experimental results.

2 Hyperreal Numbers

Hyperreal numbers were introduced in calculus to capture “infinitesimal” quantities which are infinitely small and yet not equal to zero. Formally, a number ϵ is said to be *infinitely small* or *infinitesimal* (cf. [7, 8]) iff $-a < \epsilon < a$ for every positive *real* number a . Hyperreal numbers contain all the real numbers and also all the infinitesimal numbers. There are principles (or axioms) for hyperreal numbers (cf. [8]) of which we mention:

Extension Principle.

1. The real numbers form a subset of the hyperreal numbers, and the order relation $x < y$ for the real numbers is a subset of the order relation for the hyperreal numbers.
2. There exists a hyperreal number that is greater than zero but less than every positive real number.
3. For every real function f , we are given a corresponding hyperreal function f^* which is called the *natural extension* of f .

Transfer Principle. Every real statement that holds for one or more particular real functions holds for the hyperreal natural extensions of these functions.

In short, the Extension Principle gives the *hyperreal* numbers and the Transfer Principle enables carrying out computation on them. The Extension Principle says that there does exist an infinitesimal number, for example ϵ . Other examples of hyperreals numbers, created using ϵ , are: ϵ^3 , $100\epsilon^2 + 51\epsilon$, $\epsilon/300$.

For $a, b, r, s \in \mathbb{R}^+$ and $r < s$, we have $a\epsilon^r < b\epsilon^s$, regardless of the relationship between a and b .

If $a\epsilon^r$ and $b\epsilon^s$ are used for example to denote two preference weights, then $a\epsilon^r$ is “infinitely better” than $b\epsilon^s$ even though a might be much bigger than b , i.e. coefficients a and b are insignificant when the powers of ϵ are different. On the other hand, when comparing two preferential weights of the same power, as for example $a\epsilon^r$ and $b\epsilon^r$, the magnitudes of coefficients a and b become important. Namely, $a\epsilon^r \leq b\epsilon^r$ ($a\epsilon^r > b\epsilon^r$) iff $a \leq b$ ($a > b$).

3 Keyword Preferences

We propose a framework where the user can preferentially annotate the keywords by *hyperreal numbers*.

Using hyperreal annotations is essential for reasoning in terms of “infinitely more important,” which is crucially needed in a scenario with numerous documents. This is because preference specification using only real numbers suffers from the possibility of producing senseless results as those preferences can get easily absorbed by other measures used by search engines. For instance, continuing the example given in the Introduction,

music-information-retrieval, google-search, google-ranking,

suppose that the user, dismayed of the poor result from Google, containing only three low quality pages, changes the query into¹

music-information-retrieval OR google-search OR google-ranking.

It is interesting to observe that if the user specifies this (modified) query in Google, then what she gets is a list of *many* web-pages (documents)! These pages are ranked by their Google-computed importance which is by far biased toward general pages about “google-search” and “google-ranking” rather than “music-information-retrieval.” The true pages about “music-information-retrieval” are simply buried under tons of other pages about “google-search” and “google-ranking” that are highly ranked, but contain “music-information-retrieval” either incidentally or not at all. Unfortunately, in Google and other search engines, the user cannot explicitly specify her real preferences among the specified keywords. In this example, what the user needs is a mechanism for saying that “music-information-retrieval” is of primary importance or infinitely more important than “google-search” and “google-ranking.”

But, let us suppose for a moment that Google would allow users to specify preferences expressed by real numbers. Now, imagine the user who is trying to convey that her “first and foremost” preference is for documents on “music-information-retrieval” rather than general documents about Google technology. For this, the user specifies that *music-information-retrieval* is 100 times more important than *google-search*. After all, “100 times more important” seems quite convincing in colloquial talking! However, what would happen if, according to the score computed by the search engine, general documents about *google-search* were in fact 1000 times more important than documents about *music-information-retrieval*? If the user preference levels were used to simply boost the computed document score by the same factor, then still, documents about *google-search* would be ranked higher than documents about *music-information-retrieval*. What the user would experience in this case is an “indifferent” search engine with respect to her preferences.

The solution we propose is to use hyperreal numbers for expressing preferential weights. In order to always have an effective comparison of documents with respect to a user query, we will fix an infinitesimal number, say ϵ , and build expressions on it. By the Extension Principle, such a number does exist. Now, we give the following definition.

¹ This second query style corresponds more closely than the first to what is known in the folklore as the popular “free text query:” a query in which the terms of the query are typed freeform into the search interface (cf. [6, 9]).

An *annotated free text query* is simply a set of keywords (terms) with preference weights which are polynomials of ϵ .

For all our practical purposes it suffices to consider only polynomials with coefficients in \mathbb{R}^+ . For example, $3 + 2\epsilon + 4\epsilon^2$.

By making this restriction we are able to perform symbolic (algorithmic) computations on expressions using ϵ . All such expressions translate into operations on polynomials with real coefficients for which efficient algorithms are known (we will namely need to perform polynomial additions, multiplications and divisions²).

Let us illustrate our annotated queries by continuing the above example. The user can now give

$$\textit{music-information-retrieval, google-search} : \epsilon, \textit{google-ranking} : \epsilon^2$$

to express that she wants to find documents on Music Information Retrieval and she is interested in the Google technology for retrieving and ranking music. However, by leaving intact the *music-information-retrieval* and annotating *google-search* by ϵ and *google-ranking* by ϵ^2 , the user makes her intention explicit that a document on *music-information-retrieval* is infinitely more important than any document on simply *google-search* or *google-ranking*. Furthermore, in accord with the above user expression, documents on *music-information-retrieval* and/or *google-search* are infinitely more important than documents on simply *google-ranking*. Of course, among documents on Music Information Retrieval, those which are relevant to Google search and Google ranking are more important.

We note that our framework also allows the user to specify “soft” preference levels. For example, suppose that the user changes her mind and prefers to have both *google-search* and *google-ranking* in the same “hard” preference level as determined by the power of infinitesimal ϵ . However, she still prefers, say “twice more,” *google-search* over *google-ranking*. In this case, the user gives

$$\textit{music-information-retrieval, google-search} : 2\epsilon, \textit{google-ranking} : \epsilon.$$

4 Preferentially Annotated XML Schemas

In this section, we consider the problem of weighting the structural elements of documents in a corpus with the purpose of influencing an information retrieval system to take into account the importance of different elements during the process of document ranking. Due to the wide spread of XML as a standard for representing documents, we consider in this paper XML documents which conform to a given schema (DTD). In the same spirit as in the previous section, we will use hyperreal weights to denote the importance of different elements in the schema and documents.

While the idea of weighting the document elements is old and by now part of the folklore (cf. [9]), to the best of our knowledge, there is no work that systematically studies the problem of weighting XML elements. The problem becomes challenging when elements can possibly be nested inside other elements which can be weighted as well, and one wants to achieve a consistent weight normalization reflecting the true preferences of a system administrator. Another challenging problem, as we explain in Subsection 4.4, is determining the right mapping of weights from the elements of a DTD schema into the elements of XML documents.

4.1 Hyperreal weights

In our framework, the system administrator is enabled to set the importance of various XML elements/sections in a DTD schema. For example, she can specify that the *keywords* elements of documents in an XML corpus, with “research activities” as the main theme, is more important than a section, say on *related work*. Intuitively, an occurrence of a search term in the *keywords* section is way more important than an occurrence in the *related work* section as the occurrence in the latter might be completely incidental or only loosely related to the main thrust of the document.

Thus, in our framework, we allow the annotation of XML elements by weights being, as in the previous section, polynomials of a (fixed) infinitesimal ϵ .

² The division is performed by first factoring the highest power of ϵ . For example, $(6 + 3\epsilon + 3\epsilon^2)/(4 + 2\epsilon + 3\epsilon^2)$ is first transformed into $(6\epsilon^{-2} + 4\epsilon^{-1} + 3)/(3\epsilon^{-2} + 2\epsilon^{-1} + 4)$, and then we perform the division as we would do for $(6x^2 + 4x + 3)/(3x^2 + 2x + 4)$. Observe that, as ϵ is infinitely small, ϵ^{-1} is *infinitely big*.

4.2 DTDs

Let Σ be the (finite) tag alphabet of a given XML collection, i.e. each tag is an element of Σ . Then, a DTD D is a pair (d, r) where d is a function mapping Σ -symbols to regular expressions on Σ and r is the root symbol (cf. [3]).

A *valid* XML document complying to a DTD $D = (d, s)$ can be viewed as a tree, whose root is labeled by r and every node labeled, say by a , has a sequence of children whose label concatenation, say $bc\dots x$, is in $L(d(a))$.

A simple example of a DTD defining the structure of some XML research documents is the following:

$$\begin{aligned} \text{paper} &\rightarrow \text{preamble body} \\ \text{preamble} &\rightarrow \text{title author}^+ \text{abstract keywords} \\ \text{body} &\rightarrow \text{introduction section}^* \text{related-work? references} \end{aligned}$$

where ‘+’ implies “one or more,” ‘*’ implies “zero or more” and ‘?’ implies “zero or one” occurrences of an element.

In essence, a DTD D is an extended context-free grammar, and a valid XML document with respect to D is a parse tree for D .

4.3 Annotated DTDs

To illustrate annotated DTDs, let us suppose that the system administrator wants to express that in the *body* element, the *introduction* is twice more important than a *section*, and both are infinitely more important than *related-work* and *references*, with the latter being infinitely less important than the former, we would annotate the rule for *body* as follows: $\text{body} \rightarrow (\text{introduction} : 2) (\text{section} : 1)^* (\text{related-work} : \epsilon)? (\text{references} : \epsilon^2)$.

Further annotations, expressing for example that the *preamble* element is three times more important than the *body* element, and in the *preamble*, the *keywords* element is 5 times more important than *title* and 10 times more important than the rest, would lead to having the following annotated DTD:

$$\begin{aligned} \text{paper} &\rightarrow (\text{preamble} : 3) (\text{body} : 1) \\ \text{preamble} &\rightarrow (\text{title} : 2) (\text{author} : 1)^+ (\text{abstract} : 1) (\text{keywords} : 10) \\ \text{body} &\rightarrow (\text{introduction} : 2) (\text{section} : 1)^* (\text{related-work} : \epsilon)? (\text{references} : \epsilon^2). \end{aligned}$$

Since an annotated element can be nested inside other elements, which can be annotated as well, the natural question that now arises is: How to compute the actual weight of an element in a DTD? One might be tempted to think that the actual weight of an element should be obtained by multiplying its (annotation) weight by the weights of all its ancestors. However by doing that, we could get strange results as for example a possibly increasing importance weight as we go deep down in the XML element hierarchy.

What we want here is “an element to never be more important than its parent.” For this, we propose normalizing the importance weights assigned to DTD elements. There are two ways for doing this. Either divide the weights of a rule by the sum of the rule’s weights, or divide them by the maximum weight of the rule. In the first way, the weight of the parent will be divided among the children. On the other hand, in the second way, the weight of the most important child will be equal to the weight of the parent.

The drawback of the first approach is that the more children there are, the lesser their weight is. Thus, we opt for the second way of weight normalization as it better corresponds to the intuition that nesting in XML documents is for adding structure to text rather than hierarchically dividing the importance of elements.

For example, in the above DTD, for the children of *preamble*, we normalize dividing by the greatest weight of the rule, which is 10. Normalizing in this way the weights of all the rules, we get

$$\begin{aligned} \text{paper} &\rightarrow (\text{preamble} : 1) (\text{body} : 1/3) \\ \text{preamble} &\rightarrow (\text{title} : 1/5) (\text{author} : 1/10)^+ (\text{abstract} : 1/10) (\text{keywords} : 1) \\ \text{body} &\rightarrow (\text{introduction} : 1) (\text{section} : 1/2)^* (\text{related-work} : \epsilon/2)? (\text{references} : \epsilon^2/2). \end{aligned}$$

After such normalization, for determining the actual weight of an element, we multiply its DTD weight by the weights of all its ancestors. For example, the weight of a *section* element is $(1/3) \cdot (1/2)$.

As mentioned earlier, under this weighting scheme, the most important child of a parent has the same importance as the parent itself. Thus, for instance, element *introduction* has the same importance $(1/3)$ as its parent *body*. Note that the weight normalization can of course be automatically done by the system, while we annotate using numbers that are more comfortable to write.

4.4 Weighting Elements of XML Documents

In the previous section, we described how to compute the weight of an element in a DTD. However, the weight of an element in an XML document depends not only on the DTD, but also on the particular structure of the document. This is because the same element might occur differently nested in different valid XML documents. For example, if we had an additional rule, $\text{section} \rightarrow (\text{title} : 1) (\text{text} : 1/2)$, in our annotated DTD, then, given a valid XML document, the weight of a *title* element depends on the particular nesting of this element. Namely, if the nesting is

$$\langle \text{paper} \rangle \langle \text{preamble} \rangle \langle \text{title} \rangle \dots \langle / \text{title} \rangle \dots \langle / \text{preamble} \rangle \dots \langle / \text{paper} \rangle$$

then the normalized weight of the *title* element is $1/5$. On the other hand, if the nesting is

$$\langle \text{paper} \rangle \dots \langle \text{body} \rangle \langle \text{section} \rangle \langle \text{title} \rangle \dots \langle / \text{title} \rangle \dots \langle / \text{section} \rangle \dots \langle / \text{body} \rangle \langle / \text{paper} \rangle$$

then the normalized weight of the *title* element is $(1/3) \cdot (1/2) \cdot 1 = 1/6$.

In general, in order to derive the correct weight of an element in an XML document, we need to first build the element tree of the document. This will be a parse tree for the context-free grammar corresponding to the DTD. For each node a of this tree with children $bc \dots x$, there is a unique rule $a \rightarrow r$ in the DTD such that word $bc \dots x$ is in $L(r)$.

Naturally, we want to assign weights to a 's children b, c, \dots, x based on the weights in annotated expression r . Thus, the question becomes how to map the weights assigned to the symbols of r to the symbols of word $bc \dots x$.

Since b, c, \dots, x occur in r , this might seem as a straightforward matter. However, there is subtlety here arising from the possibility of ambiguity in the regular expression. For example, suppose the (annotated) expression r is $(b : 1 + c : 1)^*(b : 2)(b : 3)^*$, and element a has three children labeled by b . Surely, bbb is in $L(r)$, but what label should we assign to each of b 's? There are three different ways of assigning weights to these b 's: $(b : 1)(b : 1)(b : 2)$, $(b : 1)(b : 2)(b : 3)$, and $(b : 2)(b : 3)(b : 3)$.

However, according to the SGML standard (cf. [4]), the only allowed regular expressions in the DTD rules are those for which we can uniquely determine the correspondence between the symbols of an input word and the symbols of the regular expression. These expressions are called "1-unambiguous" in [4].

For such an expression r , given a word $bc \dots x$ in $L(r)$, there is a unique mapping of word symbols b, c, \dots, x to expression symbols. Thus, when r is annotated with symbol weights, we can uniquely determine the weights for each of the b, c, \dots, x word symbols.

Based on all the above, we can state the following theorem.

Theorem 1. *If T is a valid XML tree with respect to an annotated DTD D , then based on the weight annotations of D , there is a unique weight assignment to each node of T .*

Now, given an XML document, since there is unique path from the root of an XML document to a particular element, we have that

Corollary 1. *Each element of a valid XML document is assigned a unique weight.*

The unique weight of an element is obtained by multiplying its local node weight with the weights of the ancestor nodes on the unique path connecting the element with the document root.³

³ All weights are considered being normalized.

5 Preferential Term Weighting and Document Scoring

Early scoring schemes were based on the Boolean model in which only the mere occurrence of terms in documents really matters. The next step was to consider the intuition that a document with more occurrences of a query term is more relevant to the query. The most popular measure reflecting this intuition is the *term frequency* (TF), which is computed as the normalized frequency of a term occurring in a document.

Formally, let V (vocabulary) be the set of distinctive terms in a collection C of documents. Denote by m and n the cardinalities of V and C respectively. Let t_i be term in V and d_j a document in C . Suppose that t_i occurs f_{ij} times in d_j . Then, the normalized term frequency of t_i in d_j is

$$tf_{ij} = \frac{f_{ij}}{\max\{f_{1j}, \dots, f_{mj}\}},$$

where the maximum is in fact computed over the terms that appear in document d_j .

Considering now XML documents whose elements are weighted based on annotated DTDs, we have that *not* all occurrences of a term “are created equal.” For instance, continuing the example in Section 4, an occurrence of a term t_i in the *keywords* element of a document is 5 times more important than an occurrence (of t_i) in the *title*, and infinitely more important than an occurrence in the *related-work* element.

Hence, we refine the *TF* measure to take the importance of XML elements into account. When an XML document conforms to an annotated DTD, each element e_k will be accordingly weighted, say by w_k .

Suppose that term t_i occurs f_{ijk} times in element e_k of document d_j . Now, we define the normalized term frequency of t_i in d_j as

$$tf_{ij} = \frac{\sum_k w_k f_{ijk}}{\max\{\sum_k w_k f_{1jk}, \dots, \sum_k w_k f_{mjk}\}}.$$

For example, suppose that t_i occurs

- once in the *keywords* element,
- twice in the *abstract* element,
- three times in the *section* elements,
- four times in the *related-work* element, and
- twice in the *references* element

of document d_j . Then, the numerator of the tf_{ij} fraction will be

$$\begin{aligned} & 1 \cdot 1 \cdot 1 + 1 \cdot (1/10) \cdot 2 + (1/3) \cdot (1/2) \cdot 3 + (1/3) \cdot (\epsilon/2) \cdot 4 + (1/3) \cdot (\epsilon^2/2) \cdot 2 \\ & = 1.7 + (2/3) \cdot \epsilon + (1/3) \cdot \epsilon^2. \end{aligned}$$

The other popular measure used in Information Retrieval is the *inverse document frequency* (IDF) which is used jointly with the TF measure. IDF is based on the fraction of documents which contain a query term. The intuition behind IDF is that a query term that occurs in numerous documents is not a good discriminator, or does not bear to much information, and thus, should be given a smaller weight than other terms occurring in few documents. The weighting scheme known as TF*IDF, which multiplies the TF measure by the IDF measure, has proved to be a powerful heuristic for document ranking, making it the most popular weighting scheme in Information Retrieval (cf. [12, 6, 9]).

Formally, suppose that term t_i occurs n_i times in a collection of n elements. Then, the *inverse document frequency* of t_i is defined to be

$$idf_i = \log \frac{n}{n_i}.$$

IDF has a natural explanation from an information theoretic point of view. If we consider a term t_i as a “message” and $p_i = \frac{n_i}{n}$ as the probability of receiving message t_i , then, in Shannon’s information theory [11], the information that the message carries is quantified by

$$I_i = -\log p_i,$$

which coincides with the IDF measure. The connection is clear; terms occurring in too many documents do not carry too much information for “discriminating” documents ([2]). On the other hand, terms that occur in few documents carry more information and hence have more discriminative power.

In XML Information Retrieval, considering each XML element that contains text as a mini-document, we can compute multiple IDF scores for a given term. Note that here, we restrict ourselves to *textual* elements only, i.e. those elements that contain terms. For instance, in the above example, *introduction* is a textual element, while *body* is not.

Depending on the importance weight of each textual element, the IDF scores should be appropriately weighted. Intuitively, in the above example, the IDF score of a term with respect to the *related-work* elements is infinitely less important than the IDF score of the term with respect to say *introduction* elements.

Formally, let E be the set of textual element-weight pairs (e_h, w_h) extracted from XML document collection C . This set is finite because C is finite, and for each element in an XML document, there is a unique weight assigned to it (see Corollary 1).

For a textual element-weight pair (e_h, w_h) , let n_h be the total number of such elements in the XML documents in collection C . Suppose that a term t_i occurs in n_{hi} of these e_h elements (of weight w_h). Then, we define the IDF of t_i with respect to these elements as

$$idf_{hi} = \log \frac{n_h}{n_{hi}}.$$

Next, we define the IDF score of a term t_i with respect to the whole document collection as

$$idf_i = \frac{\sum_h w_h \cdot idf_{hi}}{\sum_h w_h}.$$

This is the weighted average of IDF scores computed for each textual element-weight pair (e_h, w_h) .

Finally, the TF*IDF weighting scheme combines the term frequency and inverse document frequency, producing a composite weight for each term in each document. Namely, the TF*IDF weighting scheme assigns to term t_i a weight in document d_j given by

$$tf\ idf_{ij} = tf_{ij} \times idf_i.$$

In the vector space model, every document is represented by a vector of weights which are the TF*IDF scores of the terms in the document. For the other terms in vocabulary V that do not occur in a document, we have a weight of zero.

Similarly, a query q can be represented as a vector of weights with non-zero weights for the terms appearing in the query. The weights are exactly those hyperreal numbers specified by the user multiplied by the IDF scores of the terms.

Now, we want to rank the documents by computing their similarity score with respect to a query q . The most popular similarity measure is the *cosine similarity*, which for a document d_j with weight vector \mathbf{w}_j and a query q with weight vector \mathbf{w}_q is

$$\text{cosine}(\mathbf{w}_j, \mathbf{w}_q) = \frac{\langle \mathbf{w}_j, \mathbf{w}_q \rangle}{\|\mathbf{w}_j\| \times \|\mathbf{w}_q\|} = \frac{\sum_{i=1}^m w_{ij} \times w_{iq}}{\sqrt{\sum_{i=1}^m w_{ij}^2} \times \sqrt{\sum_{i=1}^m w_{iq}^2}},$$

where m is the cardinality of vocabulary V .

The above formula naturally combines the query preference weights, XML element weights, and Information Retrieval measures. Note that, we can in fact rank documents using instead the square of the cosine similarity. Thus, we only need to compare fractions of polynomial expressions based on the (fixed) infinitesimal ϵ . As such, these expressions allow for an algorithmic (symbolic) comparison procedure for ranking XML documents.

Finally, the query can be a complete document in its own. Such queries are of the type: Find all the documents which are similar to a given document. We derive weights for the elements of the query document in exactly the same manner as described in Section 4. The vector of weights for the query document is computed as for any other document in the collection. Then, this vector is compared against the vectors of the documents in the collection by computing the cosine similarity as described above.

References

1. Abiteboul S., P. Buneman and D. Suciu. *Data on the Web: From Relations to Semistructured Data and XML*. Morgan Kaufmann Publishers, San Francisco, CA, 1999.
2. Aizawa N. A. An Information-Theoretic Perspective of TF-IDF measures. *Inf. Process. Manage.* 39(1): 45–65, 2003.
3. Bex J. G., F. Neven, T. Schwentick and K. Tuyls. Inference of Concise DTDs from XML Data. *Proc. VLDB '06*, pp. 115–126.
4. Bruggemann-Klein A. and D. Wood. One-Unambiguous Regular Languages. *Inf. Comput.* 140(2): 229–253, 1998.
5. Grabs T. and H.-J. Schek. Flexible Information Retrieval on XML Documents. *Intelligent Search on XML Data*, 95–106, Springer, 2003.
6. Liu B. *Web Data Mining: Exploring Hyperlinks, Contents and Usage Data*. Springer, Berlin Heidelberg, 2007.
7. Keisler H. J. *Elementary Calculus: An Approach Using Infinitesimals*.
On-line Edition: <http://www.math.wisc.edu/~keisler/keislercalc1.pdf> 2002.
8. Keisler H. J. *Foundations of Infinitesimal Calculus*.
On-line Edition: <http://www.math.wisc.edu/~keisler/foundations.pdf> 2007.
9. Manning D. C, P. Raghavan and H. Schütze *Introduction to Information Retrieval*. Cambridge University Press. 2008.
10. McKinzie M., C. Tuckey. Higher Trigonometry, Hyperreal Numbers, and Euler’s Analysis of Infinities. *Mathematics Magazine* 74(5): 339–368, 2001.
11. Shannon C. E. A Mathematical Theory of Communication. *The Bell System Technical Journal* 27: 379-423, 1948.
12. Robertson S. Understanding Inverse Document Frequency: On theoretical arguments for IDF. *J. of Documentation* 60: 503–520, 2004.
13. Rondogiannis P., and W. W. Wadge. Minimum Model Semantics for Logic Programs with Negation-as-Failure. *ACM Trans. Comput. Log.* 6 (2): 441–467, 2005.
14. On-line Internet Shakespeare Edition. English Department, University of Victoria.
<http://internetshakespeare.uvic.ca/index.html>
15. Malik S., A. Trotman, M. Lalmas, N. Fuhr. Overview of INEX 2006. *Fifth Workshop of the INitiative for the Evaluation of XML Retrieval*, pp 1-11, 2007.

Appendix

A Experiments

Here, we describe experiments to evaluate our framework and to illustrate our ideas. For this purpose, we implemented a system incorporating our proposed framework and compared its ranking effectiveness with that of a system that ranks using the classical TF-IDF measure.

Our main research question is:

Does our preferential IR improve users' search experience compared to a traditional IR?

Here we provide practical evidence that our preferential IR does indeed perform better than a traditional IR.

As described in the previous sections, we annotated XML schema elements and search keywords in order to mark their importance in ranking the documents. We designed our experiments for both document retrieval and element retrieval. We used the following corpora as test-beds.

Corpus I On-line Internet Shakespeare Edition of the English Department ([14]), University of Victoria for element retrieval. This corpus consists of all the Shakespeare plays in XML format. The elements of interest are the speeches which total more than 33,000. For this corpus we consider all the speeches to be of the same importance, and thus, only search keyword preferences are in fact relevant for this corpus in influencing the ranking process.

Corpus II An INEX (INitiative for the Evaluation of XML retrieval) (cf. [15]) corpus. INEX is a collaborative initiative that provides reference collections (corpora). For evaluating our method, we have chosen a collection named “*topic-collection*” with numerous XML documents of moderate size. The topics of documents vary from *climate change* to *space exploration*. We preferentially annotated the DTD of this collection and gave many preferentially annotated search queries, some of which we show in this section.

A.1 Queries and Results for Corpus I

For the On-line Internet Shakespeare Edition, we created many search queries and observed that for all of them the highly ranked *speech* elements were much more relevant than the *speech* elements which were highly ranked by a traditionally implemented IR system. Here, due to space constraints, we only present two representative examples of search queries.

Q1. *romeo*, *iuliet*: ϵ , *loue*: ϵ^2 . This query says that the user is mostly interested in the keyword ‘*romeo*’ and then ‘*iuliet*’ and least interested in ‘*loue*’ (love).

Preferential IR Result for Q1. The *speech* element which was the top ranked by our preferential IR system is:

```
<s> The excellent Tragedie And Ile informe you how these things fell out. Iuliet here slaine was married to that Romeo, Without her Fathers or her Mothers grant: The Nurse was priuie to the marriage. The balefull day of this vnhappie marriage, VVas Tybalts doomesday: for which Romeo VVas banished from hence to Mantua. He gone, her Father sought by soule constraint To marrie her to Paris: but her Soule (Loathing a second Contract) did refuse To giue consent; and therefore did she vrge me Hither to finde a meanes she might auoyd What so her Father sought to force her too Or els all desperately she threatned Euen in my presence to dispatch of her selfe. Then did I giue her, (tutord my mine arte) A potion that should make her seeme as dead: And told her that I would with all post speed Send hence to Mantua for her Romeo, That he might come and take her from the Toombe, But he that had my Letters (Frier Iohn) Seeking a Brother to associate him, VVhereas the sicke infection remaind, VVas stayed by the Searchers of the Towne. But Romeo vnderstanding by his man, That Iuliet was deceasde, returnde in post Vnto Verona for to see his loue. VVhat after happened touching Paris death, Or Romeos is to me vnknowne at all. But when I came to take the Lady hence, I found them dead, and she awakt from sleep: VVhom faine I would haue taken from the tombe, VVhich she refused seeing Romeo dead. Anone I heard the watch and then I fled, VVhat after happened I am ignorant of. And if in this ought haue miscaried By of Romeo and Iuliet. By me, or by my meanes let my old life Be sacrificed some houre before his time. To the most strickest rigor of the Law. </s>
```

Traditional IR Result for Q1. The *speech* element which was the top ranked by the traditional IR system is:

<s> Consider what you first did swere vnto: To fast, to study, and to see no woman: Flat treason gainst the kingly state of youth. Say, Can you fast? your stomacks are too young: And abstinence ingenders maladies. And where that you haue vowd to studie (Lordes) In that each of you haue forsworne his Booke. Can you still dreame and poare and thereon looke. For when would you my Lord, or you, or you, Haue found the ground of Studies excellence, Without the beautie of a womans face? From womens eyes this doctrine I deriue, They are the Ground, the Bookes, the Achadems, From whence doth spring the true Promethean fire. Why vniuersall plodding poysons vp The nimble spirites in the arteries, As motion and long during action tyres The sinnowy vigour of the trauayler. Now for not looking on a womans face, You haue in that forsworne the vse of eyes: And studie too, the causer of your vow. For where is any Authour in the worlde, Teaches such beautie as a womas eye: Learning is but an adiunct to our selfe, And where we are, our Learning likewise is. Then when our selues we see in Ladies eyes, With our selves. Do we no likewise see our learning there? O we haue made a Vow to studie, Lordes, And in that Vow we haue forsworne our Bookes: For when would you (my Leedge) or you, or you? In leaden contemplation haue found out Such fierie Numbers as the prompting eyes, Of beautis tutors haue inritch you with: Other slow Artes intirely keepe the braine: And therefore finding barraine practizers, Scarce shew a haruest of their heaue toyle. But called Loues Labor's lost. But Loue first learned in a Ladies eyes, Liues not alone emured in the braine: But with the motion of all elementes, Courses as swift as thought in euery power, And giues to euery power a double power, Aboue their functions and their offices. It adds a precious seeing to the eye: A Louers eyes will gaze an Eagle blinde. A Louers eare will heare the lowest sound. When the suspitious head of theft is stopt. Loues feeling is more soft and sensible, Then are the tender hornes of Cockled Snayles. Loues tongue proues daintie, Bachus grosse in taste, For Valoure, is not Loue a Hercules? Still clymyng trees in the Hesperides. Subtit as Sphinx, as sweete and musical, As bright Appolos Lute, strung with his haire. And when Loue speaks, the voyce of all the Goddes, Make heauen drowsie with the harmonie. Neuer durst Poet touch a pen to write, Vntill his Incke were tempered with Loues sighes: O then his lines would rauish sauage eares, And plant in Tyrants milde humilitie. From womens eyes this doctrine I deriue. They sparcle still the right promethean fier, They are the Bookes, the Achademes, That shew, containe, and nourish all the worlde. Els none at all in ought proues excellent. Then fooles you were, these women to forswere: Or keeping what is sworne, you will proue fooles. For Wisedomes sake, a worde that all men loue: Or for Loues sake, a worde that loues all men. Or for Mens sake, the authour of these Women: Or Womens sake, by whom we Men are Men. Lets vs once loose our othes to find our selues, Or els we loose our selues, to keepe our othes: It is Religion to be thus forsworne. For A pleasant conceited Comedie: For Charitie it selfe fulfills the Law: And who can seuer Loue from Charitie. </s>

One can easily observe that the first speech element is clearly more relevant to the given query than the second element which is in fact quite relevant to word “loue” but not at all to the first two query keywords. We see here that the traditional TF-IDF measure has essentially ignored the first two keywords in favor of the third one just because the latter occurs too frequently in the shown document.

In the following, we show the second search query and the top-ranked speech elements for our preferential system as well as for the traditional one. For this query, similarly as for the first query, we observe that the result of the preferential system is better than that of the traditional system.

Q2. *henry*, *death*: ϵ , *king*: ϵ^2 . This query says that the user is mostly interested in the keyword ‘*henry*’ and then ‘*death*’ and least interested in ‘*king*’.

Preferential IR Result for Q2. The *speech* element which was the top ranked by our preferential IR system is:

<s> Which whiles it lasted, gaue King Henry light. O Lancaster! I feare thy ouerthrow, More then my Bodies parting with my Soule: My Loue and Feare, glew'd many Friends to thee, And now I fall. Thy tough Commixtures melts, Impairing Henry, strength'ning misproud Yorke; And whether flye the Gnats, but to the Sunne? And who shines now, but Henries Enemies? O Phoebus! had'st thou neuer giuen consent, That Phaeton should checke thy fiery Steeds, Thy burning Carre neuer had scorch'd the earth. And Henry, had'st thou sway'd as Kings should do, Or as thy Father, and his Father did, Giuing no ground vnto the house of Yorke, They neuer then had sprung like Sommer Flies: I, and ten thousand in this lucklesse Realme, Hed left no mourning Widdowes for our death, And thou this day, had'st kept thy Chaire in peace. For what doth cherrish Weeds, but gentle ayre? And what makes Robbers bold, but too much lenity? Bootlesse are Plaints, and Curelesse are my Wounds: No way to flye, no strength to hold out flight: The Foe is mercilesse, and will not pittie: For at their hands I haue deseru'd no pittie. The ayre hath got into my deadly Wounds. </s>

Traditional IR Result for Q2. The *speech* element which was the top ranked by the traditional IR system is:

<s> King. So, if a Sonne that is by his Father sent about Merchandize, doe sinfully miscarry vpon the Sea; the im- putation of his wickednesse, by your rule, should be im- posed vpon his Father that sent him: or if a Seruant, vn- der his Masters command, transporting a summe of Mo- ney, be assailed by Robbers, and dye in many irreconcil'd Iniquities; you may call the businesse of the Master the author of the Seruants damnation: but this is not so: The King is not bound to answer the particular endings of his Souldiers, the Father of his Sonne, nor the Master of his Seruant; for they purpose not their death, when they purpose their seruices. Besides, there is no King, be his Cause neuer so spotlesse, if it come to the arbitre- ment of Swords, can trye it out with all vnspotted Soul- diers: some (peradventure) haue on them the guilt of premeditated and contriued Murther; some, of begui- ling Virgins with the broken Seales of Periurie; some, making the Warres their Bulwarke, that haue before go- red the gentle Bosome of Peace with Pillage and Robbe- rie. Now, if these men haue defeated the Law, and out- runne Natiue punishment; though they can out-strip men, they haue no wings to flye from God. Warre is his Beadle, Warre is his Vengeance: so that here men are punisht, for before breach of the Kings Lawes, in now the Kings Quarrell: where they feared the death, they haue borne life away; and where they would bee safe, they perish. Then if they dye vnprouided, no more is the King guiltie of their damnation, then hee was be- fore guiltie of those Impieties, for the which they are now visited. Euery Subjects Dutie is the Kings, but euery Subjects Soule is his owne. Therefore should euery Souldier in the Warres

doe as every sicke man in his Bed, wash every Moth out of his Conscience: and dying so, Death is to him aduantage; or not dying, the time was blessedly lost, wherein such preparation was gayned: and in him that escapes, it were not sinne to thinke, that making God so free an offer, he let him out- liue that day, to see his Greatnesse, and to teach others how they should prepare. </s>

A.2 Queries and Results for Corpus II

The DTD defining the structure of this XML corpus is as follows:

```
inex_topic → title mmtitle* castitle* description narrative
```

We preferentially annotated this DTD as follows:

```
inex_topic → (title:1) (mmtitle:1/10)* (castitle:1/100)* (description:  $\epsilon$ ) (narrative:  $\epsilon^2$ ).
```

We had numerous runs on our system with preferentially annotated queries. As an example, a preferentially annotated query is as follows.

Q1. *Norway, climate: ϵ , information: ϵ^2* , where the user is looking for climate information for Norway. The query says that the user is primarily interested in keyword *Norway*, next *climate* and then, the least important, *information*.

Preferential IR Result for Q1. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
<!DOCTYPE inex_topic (View Source for full doctype...) >
- <inex_topic topic_id="447" ct_no="56" >
<title>Climate in Norway< /title>
<castitle> //article[about(., climate) and about(.,Norway)]< /castitle>
<description>Find information about the climate in Norway in summer.< /description>
<narrative>I would like to travel to Norway in july, but I have no idea about the weather. i don't know which clothes to put in my bag. To be relevant, a paragraph or a document should let me know the mean average temperature in this season and the precipitation level, or just give me an information like continental climate or polar climate...< /narrative>
< /inex_topic>
```

Traditional IR Result for Q1. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
- <inex_topic topic_id="494" ct_no="144" >
<title>ontology< /title>
<castitle> //title[about(.,ontology)]< /castitle>
<description>Find information about ontology.< /description>
<narrative>An ontology is typically a hierarchical data structure containing all the relevant entities and their relationships and rules within that domain (e.g., a domain ontology ). However, computational ontology does not have to be hierarchical at all. The computer science usage of the term ontology is derived from the much older usage of the term ontology in philosophy. For it plays a very important role in information extraction, entity recognition etc., I would like to learn more information about the introduction of it and how it works. Besides, I expect to find relevant information as elements in larger documents that deal with ontology e.g., the title of documents contains the term ontology. To be relevant, the document should contain the conception and description about ontology, something detailed about the uses of ontology as well. Information such as catalog or about specified domain without general discussion of it is not relevant.< /narrative>
< /inex_topic>
```

It is obvious that the top-ranked document of our preferential system is way more relevant than the top-ranked document of the traditional system. A similar observation applies to the other query examples which we give in the following.

Q2. *hurricane, information: ϵ* , where the user is primarily interested in keyword *hurricane*, and then *information*.

Preferential IR Result for Q2. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
- <inex_topic topic_id="530" ct_no="23" >
<title>Hurricane satellite image< /title>
<castitle> //figure[about(.,hurricane)]< /castitle>
<mmtitle> //figure[about(.,hurricane) and about(.,src:www.katrina-hurricane.biz/images/katrina-hurricane-pic3.jpg)]< /mmtitle>
<description>Find images of hurricanes taken from satellites, similar to one image from the web.< /description>
<narrative>Because I need, for a report at school on meteorological events, to have views of hurricanes taken from satellites with clues on the size of the hurricane. The images can be in greyscale or colours and we have to see the ground or at least the shape of the coasts.< /narrative>
< /inex_topic>
```

Traditional IR Result for Q2. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
- <inex_topic topic_id="494" ct_no="144">
<title>ontology< /title>
<castitle> //title[about(.,ontology)]< /castitle>
<description>Find information about ontology.< /description>
<narrative>An ontology is typically a hierarchical data structure containing all the relevant entities and their relationships and rules within that domain (e.g., a domain ontology ). However, computational ontology does not have to be hierarchical at all. The computer science usage of the term ontology is derived from the much older usage of the term ontology in philosophy. For it plays a very important role in information extraction, entity recognition etc., I would like to learn more information about the introduction of it and how it works. Besides, I expect to find relevant information as elements in larger documents that deal with ontology e.g., the title of documents contains the term ontology. To be relevant, the document should contain the conception and description about ontology, something detailed about the uses of ontology as well. Information such as catalog or about specified domain without general discussion of it is not relevant.< /narrative>
< /inex_topic>
```

Q3. *space, news*: ϵ , where the user is primarily interested in keyword *space*, and then *news*.

Preferential IR Result for Q3. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
- <inex_topic topic_id="415" ct_no="5">
<title>space history astronaut cosmonaut engineer< /title>
<castitle> //article[about(.,space history)]//section[about(., astronaut cosmonaut engineer)]< /castitle>
<description>Find the names of the 25 five most important people involved in the space exploration.< /description>
<narrative>The aim is to write a 10 pages report on the big names in the space exploration. The relevant documents should talk about at least one of the, say 25, most important people who were involved in the space exploration. Documents about one astronaut/cosmonaut who should not be personally mentioned in a 10 page report are not relevant. A relevant document should trace by itself an history of space exploration with mention of the big names, or be a document on one of these big names. So the context is space history and in this context I am looking for names of either astronauts, cosmonauts and/or engineers.< /narrative>
< /inex_topic>
```

Traditional IR Result for Q3. The top-ranked document is:

```
<?xml version="1.0" encoding="ISO-8859-1" ? >
- <inex_topic topic_id="481" ct_no="116">
<title>asia "news channel"< /title>
<castitle> //article[about(., "news channel" + asia)]< /castitle>
<description>Find articles about any of the Asian news channel.< /description>
<narrative>The TV channels which are dedicated for News alone are gaining enormous popularity. The query is aimed at finding news channels which are from asian countries. For a document to be relevant, it should include the name of the news channel along with an asian country name.If it includes more information about the news channel it will be considered more relevant.Worldwide News channels like BBC and CNN are considered as irrelevant to the query.< /narrative>
< /inex_topic>
```