

Compressing Semistructured Text Databases^{*}

Joaquín Adiego¹, Gonzalo Navarro², and Pablo de la Fuente¹

¹Departamento de Informática, Universidad de Valladolid, Valladolid, España.
{jadiago, pfuente}@infor.uva.es

²Departamento de Ciencias de la Computación, Universidad de Chile, Santiago, Chile. gnavarro@dcc.uchile.cl

Abstract We describe a compression model for semistructured documents, called *Structural Contexts Model*, which takes advantage of the context information usually implicit in the structure of the text. The idea is to use a separate semiadaptive model to compress the text that lies inside each different structure type (e.g., different XML tag). The intuition behind the idea is that the distribution of all the texts that belong to a given structure type should be similar, and different from that of other structure types. We test our idea using a word-based Huffman coding, which is the standard for compressing large natural language textual databases, and show that our compression method obtains significant improvements in compression ratios. We also analyze the possibility that storing separate models may not pay off if the distribution of different structure types is not different enough, and present a heuristic to *merge* models with the aim of minimizing the total size of the compressed database. This technique gives an additional improvement over the plain technique. The comparison against existing prototypes shows that our method is a competitive choice for compressed text databases.

Keywords: Text Compression, Compression Model, Semistructured Documents, Text Databases.

1 Introduction

The process of data compression can be split into two parts: an encoder that generates the compressed bitstream and a modeler that feeds information to it [TCB90]. These two separate tasks are called *coding* and *modeling*, respectively. Modeling assigns probabilities to symbols depending on the source data, while coding translates these probabilities into a sequence of bits. In order to work properly, the decoder must have access to the same model as the encoder.

Compression of large document collections not only reduces the amount of disk space occupied by the data, but it also decreases the overall query processing time in text retrieval systems. Improvements in processing times are achieved thanks to the reduced disk transfer times necessary to access the text in compressed form. Also, recent research on “direct” compressed text searching, i.e.,

^{*} This work was partially supported by CYTED VII.19 RIBIDI project (all authors) and Fondecyt Project 1-020831 (second author).

searching a compressed text without decompressing it, has led to a win-win situation where the compressed text takes less space and is searched faster than the plain text [WMB99,ZMNB00].

Compressed text databases pose some requirements that outrule some compression methods. The most definitive is the need for random access to the text without the possibility of decompressing it from the beginning. This outrules most adaptive compression methods such as Ziv-Lempel compression and arithmetic coding. On the other hand, semiadaptive models—which use a different model for each text encoded, building it before performing the compression and storing it in the compressed file—such as Huffman [Huf52] yield poor compression. In the case of compressing natural language texts, it has been shown that an excellent choice is to consider the words, not the characters, as the source symbols [Mof89]. Finally, the fact that the alphabet and the vocabulary of the text collections coincide permits efficient and highly sophisticated searching, both in the form of sequential searching and in the form of compressed inverted indexes over the text [WMB99,ZMNB00,NMN⁺00,MNZB00].

Although the area of natural language compressed text databases has gone a long way since the end of the eighties, it is interesting that little has been done about considering the structure of the text in this picture. Thanks to the widespread acceptance of SGML, HTML and XML as the standards for storing, exchanging and presenting documents, semistructured text databases are becoming the standard.

Our goal in this paper is to explore the possibility of considering the text structure in the context of a compressed text database. We aim at taking advantage of the structure, while still retaining all the desirable features of a word-based Huffman compression over a semiadaptive model. The idea is then to use separate semiadaptive models to compress the text that lies inside different tags.

While the possible gain due to this idea is clear, the price is that we have to store several models instead of just one. This may or may not pay off. Hence we also design a technique to *merge* the models if we can predict that this is convenient in terms of compressed file length. Although the problem of finding the optimal merging looks as a hard combinatorial problem, we design a heuristic to automatically obtain a reasonably good merging of an initially separate set of models, one per tag.

This model, which we call *Structural Contexts Model*, is general and does not depend on the coder. We plug it to a word-based Huffman coder to test it. Our experimental results show significant gains over the methods that are insensitive to the structure and over the current methods that consider the structure. At the same time, we retain all the features of the original model that makes it suitable for compressed text databases.

2 Related Work

With regard to compressing natural language texts in order to permit efficient retrieval from the collection, the most successful techniques are based on models

where the text words are taken as the source symbols [Mof89], as opposed to the traditional models where the characters are the source symbols. On the one hand, words reflect much better than characters the true entropy of the text [TCB90]. For example, a Huffman coder when words are the symbols obtains 25% versus 60% when characters are the symbols [ZMNBY00]. Another example is the WLZW algorithm (Ziv-Lempel on words) [BSTW86].

On the other hand, most information retrieval systems use words as the main information atoms, so a word-based compression eases the integration with an information retrieval system. Some examples of successful integration are [WMB99,NMN⁺00]. The text in natural language is not only made up of words. There are also punctuation, separators, and other special characters. The sequence of characters between every pair of consecutive words will be called a *separator*. In [BSTW86] they propose to create two alphabets of disjoint symbols: one for coding words and another for separators. Encoders that use this model consider texts as a strict alternation of two independent data sources and encode each one independently. Once we know that the text starts with a word or a separator, we know that after a word has been coded we can expect a separator and vice versa. This idea is known as the *separate alphabets model*.

A compression method that considers the document structure is *XMill* [LS00], developed in AT&T Labs. XMill is an XML-specific compressor designed to exchange and store XML documents, and its compression approach is not intended for directly supporting querying or updating of the compressed document. Another XML compressor is *XGrind* [TH02], which directly supports queries over the compressed files. Other approaches to compress XML data exist, based on the use of a PPM-like coder, where the context is given by the path from the root to the tree node that contains the current text. One example is *XMLPPM* [Che01], which is an adaptive compressor based on PPM, where the context is given by the structure.

3 Structural Contexts Model

Let us, for this paper, to focus on a semiadaptive Huffman coder, as it has given the best results on natural language texts. Our ideas, however, can be adapted to other encoders. Let us call *dictionary* the set of source symbols together with their assigned codes.

An encoder based on the separate alphabets model (see Section 2) must use two source symbol dictionaries: one for all the separators and the other for all the words in the texts. This idea is still suitable when we handle semistructured documents—like SGML or XML documents—, but in fact we can extend the mechanism to do better.

In most cases, natural language texts are structured in a semantically meaningful manner. This means that we can expect that, at least for some tags, the distribution of the text that appears inside a given tag differs from that of another tag. In cases where the words under one tag have little intersection with words under another tag, or their distribution is very different, the use of sep-

arate alphabets to code the different tags is likely to improve the compression ratio. On the other hand, there is a cost in the case of semiadaptive models, as we have to store several dictionaries instead of just one. In this section we assume that each tag should use a separate dictionary, and will address in the next section the way to group tags under a single dictionary.

3.1 Compressing the Text

We compress the text with a word-based Huffman [Huf52,BSTW86]. The text is seen as an alternating sequence of words and separators, where a word is a maximal sequence of alphanumeric characters and a separator is a maximal sequence of non-alphanumeric characters.

Besides, we will take into account a special case of words: *tags*. A tag is a code embedded in the text which represents the structure, format or style of the data. A tag is recognized from surrounding text by the use of delimiter characters. A common delimiter character for an XML or SGML tag are the symbols '<' and '>'. Usually two types of tags exist: *start-tags*, which are the first part of a container element, '<...>'; and *end-tags*, which are the markup that ends a container element, '</...>'.

Tags will be wholly considered (that is, including their delimiter characters) as words, and will be used to determine when to switch dictionaries at compression and decompression time.

3.2 Model Description

The structural contexts model (as the separate alphabets model) uses one dictionary to store all the separators in the texts, independently of their location. Also, it assumes that words and separators alternate, otherwise, it must insert either an empty word or an empty separator. There must be at least one word dictionary, called the *default dictionary*. The default dictionary is the one in use at the beginning of the encoding process. If only the default dictionary exists for words then the model is equivalent to the separate alphabets model.

We can have a different dictionary for each tag, or we can have separate dictionaries for some tags and use the default for the others, or in general we can have any grouping of tags under dictionaries. As explained, we will assume for now that each tag has its own dictionary and that the default is used for the text that is not under any tag.

The compression algorithm written below makes two passes over the text. In the first pass, the text is modeled and separate dictionaries are built for each tag and for the default and separators dictionary. These are based on the statistics of words under each tag, under no tag, and separators, respectively. In the second pass, the texts are compressed according to the model obtained.

At the beginning of the modeling process, words are stored in the default dictionary. When a start-structure tag appears we push the current dictionary in a stack and switch to the appropriate dictionary. When an end-structure tag

is found we must return to the previous dictionary stored in the stack. Both, start-structure and end-structure tags, are stored and coded using the current dictionary and then we switch dictionaries. Likewise, the encoding and decoding processes use the same dictionary switching technique.

3.3 Entropy Estimation

The entropy of a source is a number that only depends on its model, and is usually measured in *bits/symbol*. It is also seen as a function of the probability distribution of the source (under the model), and refers to the average amount of information of a source symbol. The entropy gives a lower bound on the size of the compressed file if the given model is used. Successful compressors get very close to the entropy.

The fundamental theorem of Shannon establishes that the entropy of a probability distribution $\{p_i\}$ is $\sum_i p_i \log_2(1/p_i)$ bits. That is, the optimum way to code symbol i is to use $\log_2(1/p_i)$ bits. In a zero-order model, the probability of a symbol is defined independently of surrounding symbols. Usually one does not know the real symbol probabilities, but rather estimate them using the raw frequencies seen in the text.

Definition 1 (Zero-order entropy estimation with multiple dictionaries)

Let N be the total number of dictionaries. The zero-order entropy for all dictionaries, \mathcal{H} , is computed as the weighted average of zero-order entropies contributed by each dictionary ($\mathcal{H}^d, d \in 1 \dots N$):

$$\mathcal{H} = \frac{\sum_{d=1}^N n^d \mathcal{H}^d}{n} \quad (1)$$

where n^d is the total number of text terms in dictionary d and n is the total number of terms that appear in the text.

4 Merging Dictionaries

Up to now we have assumed that each different tag uses its own dictionary. However, this may not be optimal because of the overhead to store the dictionaries in the compressed file. In particular, if two dictionaries happen to share many terms and to have similar probability distributions, then merging both tags under a single dictionary is likely to improve the compression ratio.

In this section we develop a general method to obtain a good grouping of tags under dictionaries. For efficiency reasons we will use the entropy as the estimation of the size of the text compressed using a dictionary, instead of actually running the Huffman algorithm and computing the exact size.

If \mathcal{V}^d is the size in bits of the vocabulary that constitutes dictionary d and \mathcal{H}^d is its estimated zero-order entropy, then the estimated size contribution of dictionary d is given by $\mathcal{T}^d = \mathcal{V}^d + n^d \mathcal{H}^d$. Considering this equation, we determine to merge dictionaries i and j when the sum of their contributions is larger than

the contribution of their union. In other words, when $\mathcal{T}^i + \mathcal{T}^j > \mathcal{T}^{i \cup j}$. To compute $\mathcal{T}^{i \cup j}$ we have to compute the union of the vocabularies and the entropy of that union. This can be done in time linear with the vocabulary sizes.

Our optimization algorithm works as follows. We start with one separate dictionary per tag, plus the default dictionary (the separators dictionary is not considered in this process). Then, we progressively merge pairs of dictionaries until no further merging promises to be advantageous. Obtaining the optimal division into groups looks as a hard combinatorial problem, but we use a heuristic which produces good results and is reasonably fast.

We start by computing \mathcal{T}^i for every dictionary i , as well as $\mathcal{T}^{i \cup j}$ for all pairs i, j of dictionaries. With that we compute the savings $\mathcal{A}^{i \cup j} = \mathcal{T}^i + \mathcal{T}^j - \mathcal{T}^{i \cup j}$ for all pairs. Then, we merge the pair of dictionaries i and j that maximizes $\mathcal{A}^{i \cup j}$, if this is positive. Then, we erase i and j and introduce $i \cup j$ in the set. This process is repeated until all the $\mathcal{A}^{i \cup j}$ values are negative.

5 Evaluation of the Model

We have developed a prototype implementing the Structural Contexts Model with a word-oriented Huffman coding, and used it to empirically analyze our model and evaluate its performance. Tests were carried out on Linux Red Hat 7.2 operating system, running on a computer with a Pentium 4 processor at 1.4 GHz and 128 Mbytes of RAM. For the experiments we selected different size collections of WSJ, ZIFF and AP, from TREC-3 [Har95].

The average speed to compress all collections is around 128 Kbytes/sec. In this value we include the time needed to model, merge dictionaries and compress. The time for merging dictionaries is included in this figure, and it ranges from 4.37 seconds for 1 Mb to 40.27 seconds for 100 Mb. The impact of merging times is large for the smallest collection (about 50% of the total time), but it becomes much less significant for the largest collection (about 5%). The reason is that it is $O(vs^2)$ to $O(vs^3)$ time, where v is the vocabulary size and s the number of different tags. Although it depends heavily on s , this number is usually small and does not grow with the collection size but depends on the DTD/schema. The vocabulary size v , on the other hand, grows sublinearly with the collection size [Hea78], typically close to $O(\sqrt{n})$.

In Figure 1 we can see a comparison for WSJ, of the compression performance using the plain separate alphabets model (SAM) and the structural context model (SCM) with and without merging dictionaries. For short texts, the vocabulary size is significant with respect to the text size, so SCM without merging pays a high price for the separate dictionaries and does not improve over SAM. As the text collection grows and the impact of the dictionaries gets reduced and we obtain nearly 11% additional compression. The SCM with merging obtains similar results for large collections (12.5% additional compression), but its performance is much better on small texts, where it starts obtaining 10.5% even for 1 Mbyte of text.

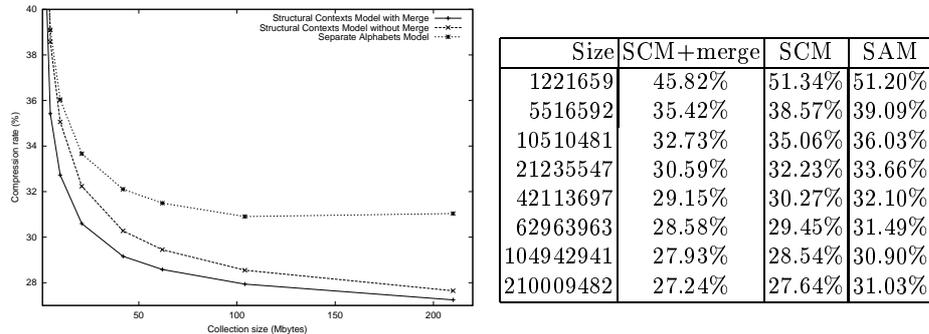


Figure 1. Compression ratios using different models, for WSJ.

Aprox.	TREC-WSJ		TREC-ZIFF		TREC-AP	
Size(Mb)	Initial	Final	Initial	Final	Initial	Final
1	11	8	10	4	9	5
5	11	8	10	4	9	5
10	11	8	10	4	9	7
20	11	9	10	6	9	7
40	11	9	10	6	9	7
60	11	9	10	6	9	7
100	11	9	10	7	9	7

Table 1. Number of dictionaries used.

Table 1 shows the number of dictionaries merged. Column “Initial” tells how many dictionaries are in the beginning: The default and separators dictionary plus one per tag, except for <DOC>, which marks the start of a document and uses the default dictionary. Column “Final” tells how many different dictionaries are left after the merge. For example, for small WSJ subsets, the tags <DOCNO> and <DOCID>, both of which contain numbers and internal references, were merged. The other group that was merged was formed by the tags <HL>, <LP> and <TEXT>, all of which contain the text of the news (headlines, summary for teletypes, and body). On the larger WSJ subsets, only the last group of three tags was merged. This shows that our intuition that similar-content tags would be merged is correct. The larger the collection, the less the impact of storing more vocabularies, and hence the fewer merges will occur. The method to predict the size of the merged dictionaries from the vocabulary distributions was quite accurate: our prediction was usually 98%–99% of the final value.

6 Conclusions and Future Work

We have proposed a new model for compressing semistructured documents based on the idea that texts under the same tags should have similar distributions. This

is enriched with a heuristic that determines a good grouping of tags so as to code each group with a separate model.

We have shown that the idea actually improves compression ratios by more than 10% with respect to the basic technique. The prototype is a basic implementation and we are working on several obvious improvements, which will make it even more competitive, especially for small collections. One is the use of canonical Huffman codes, which reduce the size of the dictionary representation. Another is a character-based compression of the vocabularies.

Other improvements would affect the results for every collection size. We can tune our method to predict the outcome of merging dictionaries: Since we know that usually our prediction is 1%–2% off, we could add a mean value to our prediction. With respect to the study of the method itself, we have to investigate more in depth the relationship between the type and density of the structuring and the improvements obtained with our method, since its success is based on a semantic assumption and it would be interesting to see how this works on other text collections.

References

- [BSTW86] J. Bentley, D. Sleator, R. Tarjan, and V. Wei. A locally adaptive data compression scheme. *Communications of the ACM*, 29:320–330, 1986.
- [Che01] J. Cheney. Compressing XML with multiplexed hierarchical PPM models. In *Proc. Data Compression Conference (DCC 2001)*, pages 163–, 2001.
- [Har95] D. Harman. Overview of the Third Text REtrieval Conference. In *Proc. Third Text REtrieval Conference (TREC-3)*, pages 1–19, 1995. NIST Special Publication 500-207.
- [Hea78] H. S. Heaps. *Information Retrieval - Computational and Theoretical Aspects*. Academic Press, 1978.
- [Huf52] D.A. Huffman. A method for the construction of minimum-redundancy codes. *Proc. Inst. Radio Engineers*, 40(9):1098–1101, 1952.
- [LS00] H. Liefke and D. Suciu. XMill: an efficient compressor for XML data. In *Proc. ACM SIGMOD 2000*, pages 153–164, 2000.
- [MNZB00] E. Silva de Moura, G. Navarro, N. Ziviani, and R. Baeza-Yates. Fast and flexible word searching on compressed text. *ACM Transactions on Information Systems*, 18(2):113–139, 2000.
- [Mof89] A. Moffat. Word-based text compression. *Software - Practice and Experience*, 19(2):185–198, 1989.
- [NMN⁺00] G. Navarro, E. Silva de Moura, M. Neubert, N. Ziviani, and R. Baeza-Yates. Adding compression to block addressing inverted indexes. *Information Retrieval*, 3(1):49–77, 2000.
- [TCB90] Ian H. Witten Timothy C. Bell, John G. Cleary. *Text Compression*. Prentice Hall, Englewood Cliffs, N.J., 1990.
- [TH02] P. Tolani and J.R. Haritsa. XGRIND: A query-friendly XML compressor. In *ICDE*, 2002. citeseer.nj.nec.com/503319.html.
- [WMB99] I.H. Witten, A. Moffat, and T.C. Bell. *Managing Gigabytes*. Morgan Kaufmann Publishers, Inc., second edition, 1999.
- [ZMNBY00] N. Ziviani, E. Moura, G. Navarro, and R. Baeza-Yates. Compression: A key for next-generation text retrieval systems. *IEEE Computer*, 33(11):37–44, November 2000.