

Statistics on The Real XML Data

Kamil Toman and Irena Mlynkova

{kamil.toman,irena.mlynkova}@mff.cuni.cz

Charles University
Faculty of Mathematics and Physics
Department of Software Engineering
Malostranske nam. 25
118 00 Prague 1, Czech Republic

Abstract. At present the eXtensible Markup Language (XML) is used almost in all spheres of human activities. We can witness a massive boom of techniques for managing, querying, updating, exchanging, or compressing XML data.

On the other hand, for majority of the XML processing techniques we can find various spots which cause worsening of their time or space efficiency. Probably the main reason is that most of them consider XML data too globally, involving all their possible features, though the real data are often much simpler. If they do restrict the input data, the restrictions are often unnatural.

We discuss the level of complexity of real XML collections and their schemes, which turns out to be surprisingly low. We involve and compare results and findings of existing papers on similar topics as well as our own analysis and we try to find the reasons for these tendencies and their consequences.

1 Introduction

Currently XML and related technologies [7] have already achieved the leading role among existing standards for data representation. They are popular for various reasons, but especially because they enable to describe the allowed structure of XML documents using powerful tools such as DTD [7] or XML Schema [9, 16, 6]. Thus we can witness a massive boom of various XML techniques for managing, processing, exchanging, querying, updating, and compressing XML data that mutually compete in speed, efficiency, and minimum space and/or memory requirements.

On the other hand, for majority of the techniques we can find various critical spots which cause worsening of their time and/or space efficiency. In the worst and unfortunately quite often case such bottlenecks negatively influence directly the most interesting features of a particular technique.

If we study the bottlenecks further, we can distinguish two typical problematic situations. Firstly, we can distinguish a group of general techniques that take into account all possible features of input XML data – an approach that is

at first glance correct. Nevertheless the standards were proposed as generally as possible enabling future users to choose what suits them most, whereas the real XML data are usually not as “rich” as they could be – they are often surprisingly simple. Thus the effort spent on every possible feature is mostly useless and it can even be harmful.

Secondly, there are techniques that somehow do restrict features of the input XML data. Hence it is natural to expect the bottlenecks to occur only in situations when given data do not correspond to the restrictions. But the problem is that such restrictions are often “unnatural”. They do not result from features of real XML data collections but from other, more down-to-earth, reasons, e.g. limitations of the basic proposal of a particular technique, complexity of such solution etc.

A solution to the given problems could be a detailed analysis of real XML data and their classification.

2 Analyses and Results

For structural analysis of XML data it is natural to view XML documents as ordered trees and DTDs or XSDs (i.e. XML Schema definitions) as sets of regular expressions over element names. Attributes are often omitted for simplicity. We use notation and definitions for XML documents and DTDs/XSDs from [8] and [5].

Up to now several papers have focused on analysis of real XML data. They analyze either the structure of DTDs, the structure of XSDs, the structure of XML data regardless their schema, or the structure of XML documents in relation to corresponding schema. The sample data usually essentially differ.

2.1 DTD vs. XML Schema

With the arrival of XML Schema, as the extension of DTD, has arisen a natural question: Which of the extra features of XML Schema not allowed in DTD are used in practise? Paper [5] is trying to answer it using analysis of 109 DTDs and 93 XSDs. Another aim of the paper is to analyze the real structural complexity for both the languages, i.e. the degree of sophistication of regular expressions used.

The former part of the paper focuses on analysis of XML Schema features. The features and their resulting percentage are:

- extension¹ (27%) and restriction (73%) of simple types,
- extension (37%) and restriction (7%) of complex types,
- **final** (7%), **abstract** (12%), and **block** (2%) attribute of complex type definitions,
- substitution groups (11%),

¹ Extension of a simple type means adding attributes to the simple type, i.e. creating a complex type with simple content.

- unordered sequences of elements (4%),
- **unique** (7%) and **key/keyref** (4%) features,
- namespaces (22%), and
- redefinition of types and groups (0%).

As it is evident, the most exploited features are restriction of simple types, extension of complex types, and namespaces. The first one reflects the lack of types in DTD, the second one confirms the naturalness of object-oriented approach (i.e. inheritance), whereas the last one probably results from mutual modular usage of XSDs. The other features are used minimally or are not used at all.

Probably the most interesting finding is, that 85% of XSDs define so called *local tree languages* [14], i.e. languages that can be defined by DTDs as well, and thus that the expressiveness beyond local tree grammars is needed rarely.

2.2 XML Document Analysis

Previously mentioned analyses focused on descriptions of the allowed structure of XML documents. By contrast paper [12] (and its extension [2]) analyzes directly the structure of their instances, i.e. XML documents, regardless eventually existing DTDs or XSDs.² It analyzes about 200 000 XML documents publicly available on the Web, whereas the statistics are divided into two groups – statistics about the Web and statistics about the XML documents.

The Web statistics involve:

- clustering of the source web sites by zones consisting of Internet domains (e.g. *.com*, *.edu*, *.net* etc.) and geographical regions (e.g. Asia, EU etc.),
- the number and volume (i.e. the sum of sizes) of documents per zone,
- the number of DTD (48%) and XSD (0.09%) references,
- the number of namespace references (40%),
- distribution of files by extension (e.g. *.rdf*, *.rss*, *.wml*, *.xml* etc.), and
- distribution of document *out-degree*, i.e. the number of **href**, **xmlhref**, and **xlink:href** attributes.

Obviously most of them describe the structure of the XML Web and categories of the source XML documents.

Statistics about the structure of XML documents involve:

- the size of XML documents (in bytes),
- the amount of markup, i.e. the amount of element and attribute nodes versus the amount of text nodes and the size of text content versus the size of the structural part,
- the amount of mixed content elements,

² The paper just considers whether the document does or does not reference a DTD or an XSD.

- the depth of XML documents and the distribution of node types (i.e. element, attribute, or text nodes) per level,
- element and attribute fan-out
- the number of distinct strings, and
- recursion.

The most interesting findings of the research are as follows:

- The size of documents varies from 10B to 500kB; the average size is 4,6kB.
- For documents up to 4kB the number of element nodes is about 50%, the number of attribute nodes about 30%. Surprisingly, for larger documents the number of attribute nodes rises to 50%, whereas the number of element nodes declines to 38%. The structural information still dominates the size of documents.
- Although there are only 5% of all elements with mixed content, they were found in 72% of documents.
- Documents are relatively shallow – 99% of documents have fewer than 8 levels, whereas the average depth is 4.
- The average element fan-out for the first three levels is 9, 6, and 0.2; the average attribute fan-out for the first four levels is 0.09, 1, 1.5, and 0.5. Surprisingly, 18% of all elements have no attributes at all.

A great attention is given to recursion which seems to be an important aspect of XML data. The authors mention the following findings:

- 15% of all XML documents contain recursive elements.
- Only 260 distinct recursive elements were found. In 98% of recursive documents there is only one recursive element used.
- 95% of recursive documents do not refer to any DTD or XSD.
- Most elements in ed pairs have the distance up to 5.
- The most common average fan-outs are 1 (60%) and 2 (37%), the average recursive fan-out is 2.2.

Last mentioned paper [11] that focuses on analysis of XML documents consists of two parts – a discussion of different techniques for XML processing and an analysis of real XML documents. The sample data consists of 601 XHTML web pages, 3 documents in DocBook format³, an XML version of Shakespeare’s plays⁴ (i.e. 37 XML documents with the same simple DTD) and documents from *XML Data repository* project⁵. The analyzed properties are the maximum depth, the average depth, the number of simple paths, and the number of unique simple paths; the results are similar to previous cases.

³ <http://www.docbook.org/>

⁴ <http://www.ibiblio.org/xml/examples/shakespeare/>

⁵ <http://www.cs.washington.edu/research/xmldatasets/>

2.3 XML Documents vs. XML schemes

Paper [13] takes up work initiated in the previously mentioned articles. It enhances the preceding analyses and defines several new constructs for describing the structure of XML data (e.g. DNA or relational patterns). It analyzes XML documents and their DTDs or XSDs eventually that were collected semi-automatically with interference of human operator. The reason is that automatic crawling of XML documents generates a set of documents that are unnatural and often contain only trivial data which cause misleading results. The collected data consist of about 16 500 XML documents of more than 20GB in size, whereas only 7.4% have neither DTD nor XSD. Such low ratio is probably caused by the semi-automatic gathering.

The data were first divided into following six categories:

- *data-centric documents*, i.e. documents designed for database processing (e.g. database exports, lists of employees etc.),
- *document-centric documents*, i.e. documents which were designed for human reading (e.g. Shakespeare’s plays, XHTML [1] documents etc.)
- *documents for data exchange* (e.g. medical information on patients etc.),
- *reports*, i.e. overviews or summaries of data (usually of database type),
- *research documents*, i.e. documents which contain special (scientific or technical) structures (e.g. protein sequences, DNA/RNA structures etc.), and
- *semantic web documents*, i.e. RDF [4] documents.

The statistics described in the paper are also divided into several categories. They were computed for each category and if possible also for both XML documents and XML schemes and the results were compared. The categories are as follows:

- *global statistics*, i.e. overall properties of XML data (e.g. number of elements of various types such as empty, text, mixed, recursive etc., number of attributes, text length in document, paths and depths etc.),
- *level statistics*, i.e. distribution of elements, attributes, text nodes, and mixed contents per each level,
- *fan-out statistics*, i.e. distribution of branching per each level,
- *recursive statistics*, i.e. types and complexity of recursion (e.g. exploitation rates, depth, branching, distance of ed-pairs etc.),
- *mixed-content statistics*, i.e. types and complexity of mixed contents (e.g. depth, percentage of simple mixed contents etc.),
- *DNA statistics*, i.e. statistics focussing on DNA patterns (e.g. number of occurrences, width, or depth), and
- *relational statistics*, i.e. statistics focussing on both relational and shallow relational patterns (e.g. number of occurrences, width, or fan-out).

Most interesting findings and conclusions for all categories of statistics are as follows:

- The amount of tagging usually dominates the size of document.

- The lowest usage of mixed-content (0.2%) and empty (26.8%) elements can be found for data-centric documents.
- The highest usage of mixed-content elements (77%) can be found for document-centric documents.
- Documents of all categories are typically shallow. (For 95% of documents the maximum depth is 13, the average depth is about 5.)
- The highest amounts of elements, attributes, text nodes, and mixed contents as well as fan-outs are always at first levels and then their number of occurrences rapidly decreases.
- Recursion is quite often, especially in document-centric (43%) and exchange (64%) documents, although the number of distinct recursive elements is typically low (for each category less than 5).
- Recursion, if used, is rather simple – the average depth, branching as well as distance of ed-pairs is always less than 10.
- The most common types of recursion are linear (20% for document-centric and 33% for exchange documents) and pure (19% for document-centric and 23% for exchange documents). On the other hand almost all schemes specify mostly general type of recursion.
- The percentage of simple mixed contents is relatively high (e.g. 79% for document-centric or even 99% for exchange documents) and thus the depth of mixed contents is generally low (on the average again less than 10).
- The number of occurrences of DNA patterns is quite high, especially for research, document-centric, and exchange documents. On the other hand the average depth and width is always low (less than 7).
- The number of occurrences of relational patterns is quite high, especially for semantic-web, research, and exchange documents. The complexity (i.e. depth and width) is again quite low.
- XML schemes usually provide too general information, whereas the instance documents are much specific and simpler.

2.4 Discussion

The previous overview of existing analyses of XML data brings various interesting information. In general, we can observe that the real complexity of both XML documents and their schemes is amazingly low.

Probably the most surprising findings are that recursive and mixed-content elements are not as unimportant as they are usually considered to be. Their proportional representation is more than significant and in addition their complexity is quite low. Unfortunately, effective processing of both the aspects is often omitted with reference to their irrelevancy. Apparently, the reasoning is false whereas the truth is probably related to difficulties connected with their processing.

Another important discovery is that the usual depth of XML documents is small, the average number is always less than 10. This observation is already widely exploited in techniques which represent XML documents as a set of points in multidimensional space and store them in corresponding data structures, e.g.

R-trees [?], UB-trees [3], BUB-trees [10] etc. The effectiveness of these techniques is closely related to the maximum depth of XML documents or maximum number of their simple paths. Both of the values should be of course as small as possible.

Next considerable fact is that the usage of schemes for expressing the allowed structure of XML documents is not as frequent as it is expected to be. The situation is particularly wrong for XSDs which seem to appear sporadically. And even if they are used, their expressive power does not exceed the power of DTDs. The question is what is the reason for this tendency and if we can really blame purely the complexity of XML Schema. Generally, the frequent absence of schema is of course a big problem for methods which are based on its existence, e.g. schema-driven database mapping methods [15].

Concerning the XML schemes there is also another important, though not surprising finding, that XML documents often do not fully exploit the generality allowed by schema definitions. It is striking especially in case of types of recursion but the statement is valid almost generally. (Extreme cases are of course recursion that theoretically allows XML documents with infinite depth or complete subgraphs typical for document-centric XML documents.) This observation shows that although XML schemes provide lots of structural information on XML documents they can be too loose or even inaccurate.

The last mentioned analysis indicates, that there are also types of constructs (such as simple mixed contents, DNA patterns, or relational patterns etc.), that are quite common and can be easily and effectively processed using, e.g., relational databases. Hence we can expect that a method that focuses on such constructs would be much effective than the general ones.

Last but not least, we must mention the problem of both syntactic and semantic incorrectness of analyzed XML documents, DTDs, and XSDs. Authors of almost all previously mentioned papers complain of huge percentage of useless sample data – an aspect which unpleasantly complicates the analyses. A consequent question is whether we can include non-determinism and ambiguity into this set of errors or if it expresses a demand for extension of XML recommendations.

3 Conclusion

The main goal of this paper was to briefly describe, discuss, and classify papers on analyses of real XML data and particularly their results and findings. The whole overview shows that the real data show lots of regularities and pattern usages and are not as complex as they are often expected to be. Thus there exists plenty of space for improvements in XML processing based on this enhanced categorization.

Acknowledgement

This work was supported in part by the National Programme of Research (Information Society Project 1ET100300419).

References

1. *The Extensible HyperText Markup Language (Second Edition)*. W3C Recommendation, August 2002. <http://www.w3.org/TR/xhtml1/>.
2. D. Barbosa, L. Mignet, and P. Veltri. Studying the XML Web: Gathering Statistics from an XML Sample. In *World Wide Web*, pages 413–438, Hingham, MA, USA, 2005. Kluwer Academic Publishers.
3. R. Bayer. The Universal B-Tree for Multidimensional Indexing: General Concepts. In *WWCA '97, Worldwide Computing and Its Applications, International Conference*, pages 198–209, Tsukuba, Japan, 1997. Springer.
4. D. Beckett. *RDF/XML Syntax Specification (Revised)*. W3C Recommendation, February 2004. <http://www.w3.org/TR/rdf-syntax-grammar/>.
5. G. J. Bex, F. Neven, and J. Van den Bussche. DTDs versus XML Schema: a Practical Study. In *WebDB '04, Proceedings of the 7th International Workshop on the Web and Databases*, pages 79–84, New York, NY, USA, 2004. ACM Press.
6. P. V. Biron and A. Malhotra. *XML Schema Part 2: Datatypes Second Edition*. W3C Recommendation, October 2004. www.w3.org/TR/xmlschema-2/.
7. T. Bray, J. Paoli, C. M. Sperberg-McQueen, E. Maler, and F. Yergeau. *Extensible Markup Language (XML) 1.0 (Third Edition)*. W3C Recommendation, February 2004. <http://www.w3.org/TR/REC-xml/>.
8. B. Choi. What are real DTDs like? In *WebDB '02, Proceedings of the 5th International Workshop on the Web and Databases*, pages 43–48, Madison, Wisconsin, USA, 2002. ACM Press.
9. D. C. Fallside and P. Walmsley. *XML Schema Part 0: Primer Second Edition*. W3C Recommendation, October 2004. www.w3.org/TR/xmlschema-0/.
10. R. Fenk. The BUB-Tree. In *VLDB '02, Proceedings of 28th International Conference on Very Large Data Bases*, Hong Kong, China, 2002. Morgan Kaufman Publishers.
11. J. Kosek, M. Kratky, and V. Snasel. Struktura realnych XML dokumentu a metody indexovani. In *ITAT 2003 Workshop on Information Technologies Applications and Theory*, High Tatras, Slovakia, 2003. (in Czech).
12. L. Mignet, D. Barbosa, and P. Veltri. The XML Web: a First Study. In *WWW '03, Proceedings of the 12th international conference on World Wide Web, Volume 2*, pages 500–510, New York, NY, USA, 2003. ACM Press.
13. I. Mlynkova, K. Toman, and J. Pokorny. *Statistical Analysis of Real XML Data Collections*. Technical report 2006/5, Charles University, June 2006. <http://kocour.ms.mff.cuni.cz/~mlynkova/doc/tr2006-5.pdf>.
14. M. Murata, D. Lee, and M. Mani. Taxonomy of XML Schema Languages using Formal Language Theory. In *Extreme Markup Languages*, Montreal, Canada, 2001.
15. J. Shanmugasundaram, K. Tufte, C. Zhang, G. He, D. J. DeWitt, and J. F. Naughton. Relational Databases for Querying XML Documents: Limitations and Opportunities. In *VLDB'99, Proceedings of 25th International Conference on Very Large Data Bases*, pages 302–314, Edinburgh, Scotland, UK, 1999. Morgan Kaufmann.
16. H. S. Thompson, D. Beech, M. Maloney, and N. Mendelsohn. *XML Schema Part 1: Structures Second Edition*. W3C Recommendation, October 2004. www.w3.org/TR/xmlschema-1/.