

jInfer XML Schema Inference Framework

Michal Klempa, Mário Mikula, Robert Smetana, Michal Švirec, Matej Vitásek
Advisors: RNDr. Irena Mlýnková, Ph.D., Martin Nečaský, Ph.D.

Praha, 2011

Abstract

Today, many algorithms solving XML schema inference problem exist. None of them is widely used by public to practically solve the problem. This is because the lack of simple user interface to algorithms. We present NetBeans ([net]) based pluggable framework to fill this gap. Thus enabling people dealing with XML in practice to obtain schema for their documents by following simple steps, without need to read and understand theoretical aspects. We hope our solution will boost usage of XML schemas in practice, since creating schema for existing set of documents will be more affordable (one doesn't have to write one by hand).

Related work

Short architecture overview

More work/Open issues

Conclusion?

Introduction

Although many algorithms ([Aho96, BNST06, BNV07, VMP08, Vyh]) try to solve problem of inferring schema for existing set of XML documents, have you ever tried to solve the problem in practice?

When the problem arises, one wants to find solution as quickly as possible, without need to investigate many scientific papers. We focus on this group of potential users - programmers, system administrators, xml coders, xml maintainers in private and scientific sector.

Framework is designed extensible and is implemented as a set of plug-ins for NetBeans platform. This brings second group of users (although much smaller in counts) of framework - scientific developers willing to experiment with their new algorithms may easily gain benefits of user interface provided, thus speeding up their development, easing them comparison with other algorithms already developed for framework.

But it is not all about user interface. We bring in XML/XSD/DTD import/export modules ready to use. We implement sample inferring algorithm (see [Aho96]). We provide visualization tools to visualize input XML documents as a set of grammar rules, and tools to visualize non-deterministic finite automata used in many of inferring algorithms. Thus extending existing NFA based algorithm with user interaction is made simple for scientific developers.

References

- [Aho96] H. Ahonen. *Generating grammars for structured documents using grammatical inference methods*. PhD thesis, Department of Computer Science, University of Helsinki, Series of Publications A, Report A-1996-4, 1996.
- [BNST06] Geert Jan Bex, Frank Neven, Thomas Schwentick, and Karl Tuyls. Inference of concise dtlds from xml data. In *Proceedings of the 32nd international conference on Very large data bases, VLDB '06*, pages 115–126. VLDB Endowment, 2006.
- [BNV07] Geert Jan Bex, Frank Neven, and Stijn Vansummen. Inferring xml schema definitions from xml data. In *Proceedings of the 33rd international conference on Very large data bases, VLDB '07*, pages 998–1009. VLDB Endowment, 2007.
- [net] Todo. <http://netbeans.org>.
- [VMP08] Ondřej Vošta, Irena Mlýnková, and Jaroslav Pokorný. Even an ant can create an xsd. In *DASFAA'08: Proceedings of the 13th international conference on Database systems for advanced applications*, pages 35–50, Berlin, Heidelberg, 2008. Springer-Verlag.
- [Vyh] Julie Vyhnanovská. Automatic construction of an xml schema for a given set of xml documents.