

Paper 042-2008

XML for SAS® Programmers

Frederick Pratter, Computer Science and Multimedia Studies Program,
Eastern Oregon University, La Grande OR

ABSTRACT

XML (the eXtended Markup Language) is an open standard for the definition, transmission, validation, and interpretation of data. The standard was developed by the Worldwide Web Consortium (W3C) in order to provide a simple and efficient way to manage self-documenting data files. SAS® 9.1.3 includes a number of useful tools for creating and parsing XML data. These include:

- The XML libname engine - used to import and export documents in XML format into or from a SAS dataset, optionally using an XMLMap.
- The XML Mapper application - a graphical interface to analyze the structure of an XML document or an XML schema and generates basic XML syntax for the XMLMap.
- The ODS MARKUP output destination - used to create XML from SAS output; ODS MARKUP creates but does not read XML documents.

This paper will introduce these concepts and present some brief examples of each. In addition, bonus topics include using PROC CDISC, the CDISC XML formats and creating XML files for Microsoft Office using the ExcelXP ODS tagset.

INTRODUCTION

It is self-evident that the Web browser has become the standard graphical user interface for the “Naughties”. Platform independent, universal, and free, the HTML browser is the default technology for collecting and displaying information. Unfortunately, this ubiquity has resulted in a number of unavoidable consequences for those of us who have to produce the software to support these activities. For one, everybody has had to learn HTML. The big advantage is that the playing field has been leveled—we are all working with an equal handicap.

Fortunately, SAS has risen to the challenge, introducing a number of new tools, including the Output Delivery System, SAS/IntrNet®, and the AppDev Studio™, that can make the programmer’s life easier. A full review of ODS tagsets, webAF™, htmSQL, and the Application Dispatcher would be beyond the scope of what can be covered in a one hour presentation (but see *Web Development with SAS by Example, 2nd. ed.*, SAS Books by Users, 2006 for more information). SAS has done a reasonably good job of documenting the XML features. The purpose of this paper is to provide some background and to explain some of the terminology used in the documentation.

An important component of these new standards is XML, the eXtended Markup Language. This paper addresses one specific aspect of application development using SAS, focusing on the SAS utilities that are available for parsing and generating XML. In the process, the following questions will be considered:

- what is XML?
- what is the difference between XML and HTML?
- what does XML do?
- what SAS tools are available for XML processing?
- what do they do, and why would you want to use them?

The main difference between XML and HTML is that they were designed with different goals in mind. HTML was intended as a way of displaying information. As such, it includes both the data and the presentation formatting. XML was designed to separate the data from the presentation; it contains only the data. It is important to recognize that XML is not a replacement for HTML. They are used for different purposes, and consequently it is likely that the two markup languages will continue to coexist peacefully.

There are a number of papers available on using SAS with XML from past user group meetings; see the references at the end of this paper for a partial list. Also, there are at least five presentations on using SAS with XML scheduled for this conference, including “Power Up XML with SAS®” by Richard Foley and Tony Friebel in the *SAS Presents* section. The “SAS Notes and Concepts for ODS” document available at http://support.sas.com/rnd/base/ods/templateFAQ/Template_xml.html includes an extensive discussion on using XML, ODS, and the MARKUP Destination. SAS also offers a Live Web course on “Creating Markup Language Files Using ODS Markup, SAS XML LIBNAME Engine, and the TEMPLATE Procedure” (see <http://support.sas.com/training/us/crs/lwxml.html>).

XML BASICS

XML is a very simple idea really. Its primary purpose is the exchange of data between dissimilar systems. As such, it is an extension of earlier standards, such as comma delimited (CSV) files. A plain text file containing variable names on the first line and rows of data separated by commas has been a standard means for data transfer since the early pre-DOS days. The well known drawback of comma delimited files is that of course you cannot send text strings that contain commas. You can always surround the text with quotes, but then you cannot send strings containing quotes. XML makes it all so much easier, since you can send almost anything enclosed between opening and closing tags. (If you need to send the angle brackets "<" and ">", the HTML *entities* < and > can be used.)

There is another, more important advantage of XML over CSV format. This is that the former can be used to structure hierarchical data, whereas CSV files must be flat tables. (This causes a problem for SAS, of which more later.) This is an immense improvement, since it allows data to be exchanged at any level of complexity. The receiving system need know nothing whatever about the data, except that it is in XML format, since the document carries its own description with it.

The following code shows an example of self-documenting XML, loosely borrowed from Paul Kent's original SUGI 26 presentation.

```
<?xml version="1.0" encoding="ISO-8859-1"?>
<conference>
  <talk>
    <name>Paul</name>
    <title session="Advanced Tutorials">XML and SAS</title>
    <when>Tuesday, 3pm</when>
    <where>hope to find out </where>
  </talk>
  <talk>
    <name>Paul</name>
    <title session="Forums">Futures Forum</title>
    <when>Wednesday, 5 pm</when>
    <where>Room 220</where>
  </talk>
</conference>
```

This little example illustrates most of the rules of XML. The first statement is required, indicating that the document is in XML, that the version is 1.0 (actually, there is no other version), and that the character encoding is the ISO Latin-1/West European character set.

All XML documents must have a root tag; in this case it is <conference>. Unlike HTML, every XML element must also have a closing tag. (An empty tag, that is one with no associated data, can be abbreviated, as for example
 in HTML.) Thus the last line of the example is </conference>, which closes the root tag. White space is ignored and elements can begin and end anywhere on the line.

Tags can be nested, but they have to be closed in the opposite order from the opening tags. Thus while

```
<h1><b>Hello World</h1></b>
```

may be valid HTML, it is illegal in XML.

Next, XML tags are case sensitive: <message> cannot be closed with </MESSAGE>.

XML elements can have attributes, just like HTML. However in XML the attribute value must always be quoted, for example <title session="Forums">. Attributes in XML are used to convey information that describes the elements; it is not always clear what should be an element and what an attribute.

PARSING XML

There is a special set of standards for reading ("parsing") well-formed XML documents. These standards are embodied in a number of freely available programs for traversing XML documents to extract the information contained therein. Most of these are written in Java, although that is not at all necessary. It would certainly be possible to write your own XML parser in a DATA step, but the general feeling is that given the availability of well tested ones, why bother?

Parsers come in two flavors (DOM, the Document Object Model, and SAX, the Simple API for XML). The SAS® 9.1.3 XML LIBNAME Engine User's Guide (<http://support.sas.com/rnd/base/xmlengine/sxle913/usersguide913.htm>) has an FAQ section that discusses the different parsers and how they are used in SAS:

Currently, the XML engine can be either a DOM application or a SAX application, depending on what you are doing:

- If the format type is either GENERIC (the default) or ORACLE, the XML engine uses a modified Document Object Model (DOM), which converts the document's contents into a node tree. However, for the XML engine, the node tree cannot be queried (traversed).
- If you are using an XMLMap to import an XML document, the XML engine uses a Simple API for XML (SAX) model. SAX does not provide a random access lookup to the document's contents; it scans the document sequentially and presents each item to the application only one time.

Note that for large XML documents for which you are simply using the format type GENERIC or ORACLE, if you are having resource problems, convert to using an XMLMap, which uses the SAX model.

The difference between these is well beyond the scope of this paper, but a useful introduction to XML parsers can be found at www.ibiblio.org/xml/slides/gsdcd/fundamentals. XML parsers are available for downloading from Microsoft and Oracle, as well as from the Apache Software Foundation. In addition, SAS has provided a reasonably powerful set of tools for parsing and generating XML. The explanation of the XML libname engine below, explains how to parse documents in SAS.

VALIDATING XML DOCUMENTS

The rules for XML, if followed correctly, result in a document that is "well formed". An XML document that is not well formed cannot be parsed—the application should return an error if, for example, a tag is not properly closed. But XML documents can also be *validated*. A valid XML document is one that corresponds to a specified *Document Type Definition*, or DTD.

It is not necessary to have a DTD to use XML, but it adds a layer of data integrity checking that can be very useful. The purpose of a DTD is to specify the legal elements of an XML document. A DTD can be declared inline within the XML, or as an external reference. The following example shows an external DTD for the XML shown above:

```
<!DOCTYPE talk [
  <!ELEMENT talk      (name, title, when, where)>
  <!ELEMENT name      (#PCDATA)>
  <!ELEMENT title      (#PCDATA)>
  <!ELEMENT when       (#PCDATA)>
  <!ELEMENT where      (#PCDATA)>
  <!ATTLIST title      session CDATA #REQUIRED>
]>
```

The Document Type Definition indicates which elements can appear inside a legal `<talk>` element. The elements are further described as "parsed character" data, that is, as text that will be treated as markup. One of the elements, `<title>` has a required character attribute `session`. It would be possible to specify that in addition to being required, `session` must also be one of the values from an enumerated list, but this was not done in the example. Clearly, a DTD can be used to verify that the data you receive is valid, just as others can validate the data you send them.

In order to validate the sample XML document, it is necessary to save the DTD in some public location and then include in the XML document a reference to the DTD, as for example:

```
<!DOCTYPE note SYSTEM "talk.dtd">
```

This indicates to the parser that the rules in the specified DTD must be observed. This is a `SYSTEM` DTD, indicating that it is located locally; the pathname above references a file in the same directory as the XML, which is likely to be the case. In addition, a number of standard DTDs have been developed by independent groups responsible for data exchange. These are public DTDs, as for example:

```
<!DOCTYPE HTML PUBLIC "-//W3C//DTD HTML 4.01//EN"
"http://www.w3.org/TR/html4/strict.dtd">
```

This Document Type Definition specifies legal HTML elements so that an HTML document can be validated as XML.

The problem with DTDs is that they are not themselves XML, which seems peculiar. *XML Schema* is an XML-based alternative to DTDs that describes the structure of an XML document in XML. The XML Schema language is also referred to as *XML Schema Definition* (XSD). Schemas are harder to write than DTDs, so they have been slow to catch on, but they offer a richer and more powerful syntax, so they are used for applications where validation is important, as in the CDISC standard.

The XML libname engine supports DTDs and schemas on output but not, unfortunately, on input; the engine just assumes that input documents are valid and well-formed.

SAS XML LIBNAME ENGINE

The XML libname engine can be used both to import and export XML documents. That is, a SAS data set can be written out in XML, and, providing it is properly structured, an XML document can be read into a SAS dataset.

Creating an XML document is rather difficult in Java, and somewhat tricky in Oracle. In SAS it is a lead pipe cinch, as the following illustrates:

```
libname oracledb ORACLE user=scott password=tiger path='sample';

libname xmltrans XML "C:\My Documents\XML\all_users.xml";

data xmltrans.all_users;
  set oracledb.all_users;
run;
```

This simple program creates an XML document from a database table. Note that the input file in this case is an Oracle library—it could be SAS, MS/Access, or anything else that SAS offers a libname engine for.

The XML libname engine also provides support for CDISC formats. The Clinical Data Interchange Standards Consortium (CDISC) develops XML standards for the pharmaceutical industry. To export a CDISC XML document using an XML libname, the *User's Guide* offers the following code (see "Exporting an XML Document in CDISC ODM Format"):

```
filename output 'C:\myoutput.xml';
libname output XML xmltype=CDISCODM
  formatactive=yes;
data output.AE2;
  set odm.AE;
```

You can also import a CDISC XML file, as shown below (see "Importing a CDISC ODM Document"):

```
filename odm 'C:\Documents and Settings\myid\My Documents\CDISC\AE.XML';
libname odm XML xmltype=CDISCODM
  formatActive=YES formatNoReplace=NO formatLibrary="Work"
```

In addition to the XML libname engine, support for CDISC XML documents is available in PROC CDISC. Two formats are supported: the *Operational Data Model (ODM) 1.2* and the *Study Data Tabulation Model (SDTM) 1.3*; see <http://support.sas.com/rnd/base/xmlengine/proccdisc/TW8774.pdf> and <http://support.sas.com/rnd/base/xmlengine/proccdisc/cdiscsdm.html> for more information on using PROC CDISC to read and write XML.

USING XML MAPPER APPLICATION

Parsing XML in SAS is straightforward as long as the document follows the required structure (see "Understanding the Required Physical Structure for an XML Document to Be Imported Using the GENERIC Format Type"):

- The root-enclosing element (top-level node) of an XML document is the document container. For SAS, it is like the SAS library.
- The nested elements (repeating element instances) that occur within the container begin with the second-level instance tag.
- The repeating element instances must represent a rectangular organization. For a SAS data set, they determine the observation boundary that becomes a collection of rows with a constant set of columns.

If your XML document does not meet these requirements, it is necessary to turn to transform it so that it does. SAS now offers a standalone utility called the XML Mapper. As the *User's Guide* indicates:

SAS XML Mapper is a Java-based application that assists you in creating and modifying XMLMaps for use by the XML engine. SAS XML Mapper provides a graphical interface that you can use in order to generate the appropriate XML elements. SAS XML Mapper analyzes the structure of an XML document or an XML schema and generates basic XML syntax for the XMLMap. The interface consists of windows, a menu bar, and a tool bar. Using SAS XML Mapper, you can display an XML document or an XML schema, create and modify an XMLMap, and generate example SAS programs.

The problem is that by nature, SAS datasets are tabular, that is the data are in rows and columns, but XML documents may be hierarchical. The XML Mapper application is used to pick and choose the variables and observations from the XML. SAS views an XML document stored in a single file as a library, so more than one map can be constructed from a single document, at different levels of aggregation. The resulting map is specified to the engine by the `xmlmap=` option on the libname statement.

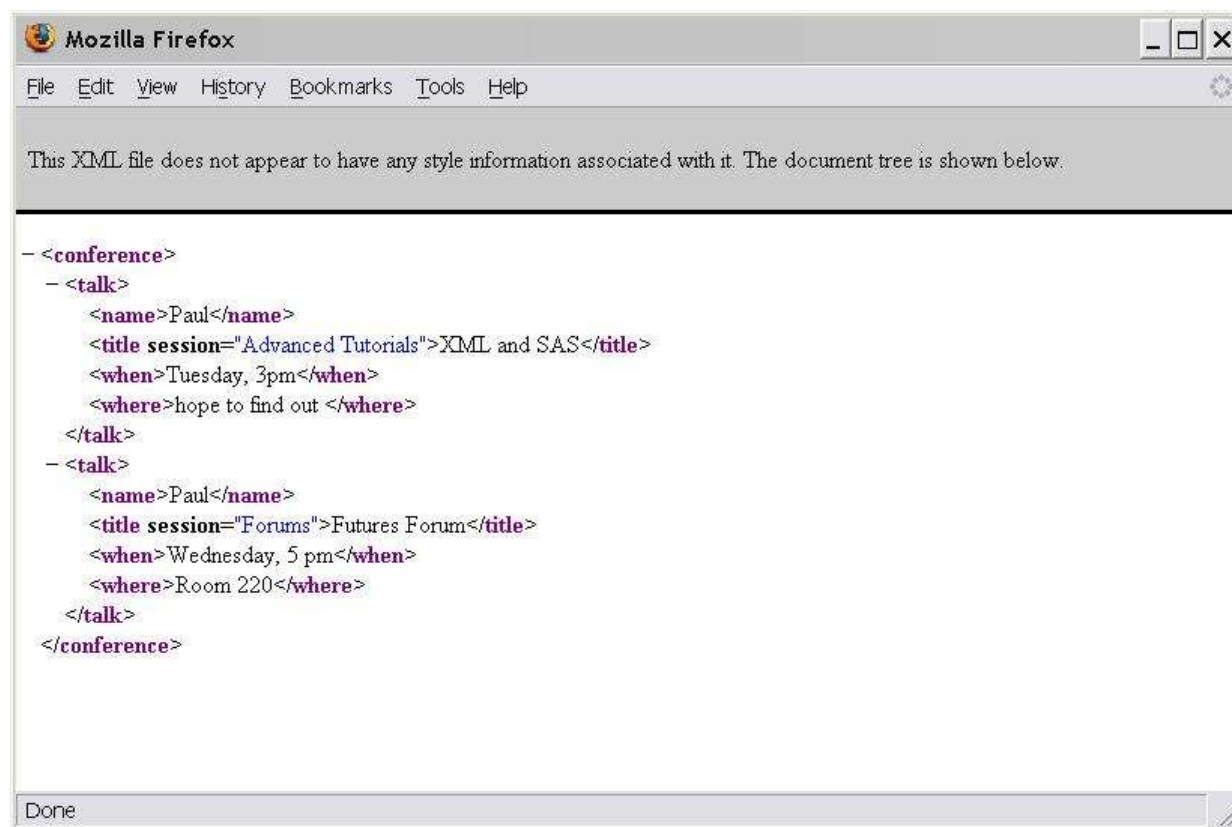
For most users, the XML engine will provide sufficient functionality as is, at least for transferring data files. In addition, it is currently the only way (other than writing a custom DATA step) to read in XML documents. For more complex applications, and in particular for formatting procedure output, the Output Delivery System has been enhanced to support a wide range of output options.

FORMATTING XML DOCUMENTS

There is one more topic that is important for an understanding of how XML can be used, one which allows the introduction of two more acronyms-- CSS and XSL. *Cascading stylesheets* (CSS) were introduced in HTML 4.0 to specify document formatting. Style tags define how HTML elements are displayed. In this they are an extension of the `` tag and `color=` attribute in earlier HTML versions.

While styles can be included in the document, it is more usual to save the styles in external files. In this way, a common "look and feel" can be maintained across a Web site. If all of the developers use a common set of style sheets, it is possible to change the appearance and format of all the pages in the site just by editing a single CSS documents. (They are called 'cascading' because if more than one stylesheet is used, each successive style definition will add new values and over-ride values previously defined.)

Note that Microsoft Internet Explorer 6 comes with a built-in style sheet that can be used to display XML as a tree structure. (It doesn't work in EE 7!) In order to be able to view XML as anything other than the browser default, an HTML LINK statement needs to be added to the XML, pointing to a custom stylesheet for that site. Note that Firefox displays a warning message that it is using the default style for displaying the page:



Although it is certainly possible to use CSS for displaying XML, there is one disadvantage, which is that CSS requires the limited set of tags that are defined in HTML. Because XML does not use predefined tags, it is necessary to have a way to indicate how the document is to be displayed. If you want your document to show up as a table it is necessary to use XSL, the eXtensible Stylesheet Language. This is a stylesheet language specifically developed for formatting XML for display.

XSL is generally conceded to be a lot of bother for very little marginal improvement over what you can do with CSS. The XSL transformation language (XSLT) on the other hand, has been widely adopted and continues to be a focus of considerable developer interest. XSLT is used, just like it sounds, to transform a XML document into something else, such as a different XML document, HTML, or even CSV. What is more, XSLT can filter and sort XML, address parts of an XML document, and output XML to different devices. The Output Delivery System supplies a powerful and

easy-to-use set of XML transformations that are far simpler to use than XSLT but which provide much of the same functionality.

XSLT represents the most sophisticated tool available for creating and manipulating markup documents. Most users will simply want to parse and or generate XML, and for that SAS has supplied simple and effective mechanisms in the XML libname engine and the ODS MARKUP output destination.

USING THE OUTPUT DELIVERY SYSTEM TO CREATE XML

SAS Release 8 introduced the ODS XML driver as an experimental driver, with no guarantee that the output would be valid XML. There was only one DTD available, which produced a document in a single standard format. With Release 9.1, SAS provided the new ODS MARKUP statement, allowing the user to export a variety of markup languages, including HTML, XML, CSV, DTD, CSS and XSL. The ODS MARKUP statement uses essentially the same syntax as the deprecated ODS HTML statement, except for the addition of TAGSET= option. The value of this option determines the type of output file to be created. (Again, it is important to note that ODS cannot be used to parse XML, only generate it.)

In addition to the list of tagsets available from SAS, it is also possible to create new ones, as well as customizing the SAS-supplied tagsets. The new TEMPLATE procedure is used to review, create, and customize tagsets. Creating an XML document using ODS is as simple as the following illustration:

```
ods listing close;
ods markup body='C:\My Documents\xml\class.xml';
proc print data=sashelp.class;run;
ods markup close;
```

In order for the recipient to validate the resulting XML, it is also possible to create the DTD at the same time as the XML, as the following example illustrates:

```
libname myfiles 'C:\My Documents\myfiles';
ods listing close;
ods markup body='C:\My Documents\xml\statepop.xml'
  frame='C:\My Documents\xml\statepop.dtd' tagset=default;
proc univariate data=myfiles.statepop;
  var citypop_90 citypop_80;
  title 'US Census of Population and Housing';
run;
ods markup close;
```

With XML it is necessary to have some kind of style sheet to be able to view this output as anything other than plain text. The obvious question arises, why go to all this trouble when ODS HTML will produce a satisfactory Web page?

The answer requires a digression into static versus dynamic HTML content. Using ODS HTML to output the results of a procedure produces a static HTML document. If you want to change the formatting it is possible to do so, but re-running the procedure results in a new document with the original format. A certain amount of customizing is possible, but the ODS HTML statement produces a default format each time it is invoked. The obvious workaround is to use a stylesheet. Each time the procedure is run, a new HTML document will be generated, but the formatting will come from the stylesheet and not the document. The same process is much more efficient in XML, however. Rather than generating a default format and then overriding it, an XML document by its nature contains no inherent display format. Depending on the complexity of the report, a CSS or XSL stylesheet could be used to produce the output result. The programmer need only rerun the SAS code to generate a new XML document every time. This turns out to be pretty straightforward, using some of the new SAS/IntrNet and AppDev studio tools, but that is a topic for another paper.

One more ODS topic that needs to be mentioned is the new ExcelXP tagset. SAS has provided a demo at http://support.sas.com/rnd/base/ods/odsmarkup/excelxp_demo.html; Vince DelGobbo will be presenting at this conference on "Creating AND Importing Multi-Sheet Excel Workbooks the Easy Way with SAS®." Base SAS software in Release 9.1 and later can be used to easily transfer data from SAS files and procedures to MS Office using ODS and the ExcelXP tagset. The SAS documentation includes a handy introduction "TEMPLATE Procedure: Creating Markup Language Tagsets" that explains how to customize and modify your own tagsets.

CONCLUSION

XML is the new industry standard for data transfer between dissimilar platforms. In addition, it is an alternative to HTML for data display, when used with custom stylesheets to create a formatted document by the client's browser. SAS has provided the XML libname engine for the first function, and has modified the Output Delivery System in order to accommodate the second. These features are new in the most recent SAS releases, and are the subject of ongoing development. This paper has reviewed the current level of functionality of these tools. As always, SAS has

solicited user input in determining the best mix of features and usability. It is important for the user community to try to begin to integrate these tools into practical applications, in order to drive the next generation of SAS XML utilities.

REFERENCES

In addition to the references included in the paper, the following conference presentations provide useful additional explanations.

- Cisternas, Miriam and Ricardo Cisternas. "Reading and Writing XML files from SAS[®]," SUGI 29 (2003), <http://www2.sas.com/proceedings/sugi29/119-29.pdf>.
- Friebel, Anthony. "XML? We do that!" SUGI 28 (2002), <http://www2.sas.com/proceedings/sugi28/173-28.pdf>.
- Gebhart, Eric S. "ODS MARKUP: The SAS[®] Reports You've Always Dreamed Of," SUGI 30 (2004), <http://www2.sas.com/proceedings/sugi30/085-30.pdf>.
- Hoyle, Larry. "Reading Microsoft Word XML files with SAS[®]," SUGI 31 (2005), <http://www2.sas.com/proceedings/sugi31/019-31.pdf>.
- Palmer, Michael. "XML in the DATA Step," SUGI 29 (2003), <http://www2.sas.com/proceedings/sugi28/025-28.pdf>.
- Pratter, Frederick. "Beyond HTML: Using the SAS System[®] Version 8.2 with XML," SUGI 27 (2001), <http://www2.sas.com/proceedings/sugi27/p002-27.pdf>.
- Pratter, Frederick. "Using the SAS[®] Output Delivery System and PROC TEMPLATE to Create XHTML Files," SAS Global Forum 2007, <http://www2.sas.com/proceedings/forum2007/118-2007.pdf>.
- Shoemaker, Jack N. "XML Primer for SAS[®] Programmers," SUGI 30 (2004), <http://www2.sas.com/proceedings/sugi30/240-30.pdf>.

CONTACT INFORMATION

Frederick Pratter
Computer Science/Multimedia Studies Program
Eastern Oregon University
One University Blvd.
La Grande OR 97850
fpratter@eou.edu
<http://cs.eou.edu/CSMM/fpratter>

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration. Other brand and product names are registered trademarks or trademarks of their respective companies.