

## XML (extensible Markup Language).

XML Stands for extensible Markup language. It is a text-based mark-up language derived from Standard Generalised Markup language (SGML). XML tags identify the data and are used to store and organise the data, rather than specifying how to display it like HTML tags, which are used to display the data.

XML can work behind the scene to simplify the creation of HTML documents for large websites. It can be used to exchange the information between organisations and systems. XML can be used to store and arrange the data and it is also used for offloading and reloading the databases. XML can easily be merged with style sheets to create almost any desired output.

### Advantages of XML :-

- 1.) XML document is human readable and we can edit any XML document in simple text editors.
- 2.) The XML document is language neutral. That means a Java program can generate an XML document.
- 3.) Every XML document has a tree structure. Hence, complex data can be arranged systematically & can be understood in simple manner.
- 4.) XML files are independent of an operating system.

## features of XML :-

- XML is intended for transport (or) storage of the data.
- The most important feature of XML is that we is able to define & use his own tag.
- XML contains only data & does not contain any formatting information.
- XML permits the document writer to create the tags for any type of information.
- Searching, Sorting, rendering (or) manipulating the document is possible using extended Stylesheet language (XSL).
- Some commonly used web browsers like internet Explorer & mozilla firefox provide support to the tags in XML.

## Key Points :-

- 1.) XML is case sensitive.  
Eg:- <message> hi </message> (WRONG).  
<Message> hi </Message> (CORRECT).
- 2.) In XML Start tag must have matching end tag.  
Eg:- <Starting> Hi </end> (WRONG)  
<Starting> Hi </Starting> (CORRECT).
- 3.) The elements in XML must be properly nested.  
Eg:- <one> <two> Hi </one> </two> (WRONG)  
<one> <two> Hi </two> </one> (CORRECT)
- 4.) XML must be saved with ".xml".

Comments in XML:-

The syntax of writing comments in XML is similar to HTML comments.

Eg: `<!-- This line is a comment -->`.

Simple XML pgm:-

`<?xml version = "1.0"?>`

`<employees>` → Root element

`<employee>` → Container element

`<name> Sai Eswari </name>`

`<dept> IT </dept>`

`</employee>`

`</employees>`

Defining XML TAGS:-

In XML the basic entity is element. The elements are used for defining the tags. Mostly only one element is used to define a single tag.

We can broadly categorize XML tags as follows.

i.) Start tag: The beginning of every non-empty XML element is marked by a start-tag.

Ex: `<element name>`

ii.) End tag: Every element that has a start tag should tag end with an end tag.

Ex: `</element name>`

iii.) Empty tag: An element which has no content is known as empty. An empty tag can be

defined by putting a '/' (forward slash) before closing bracket.

Ex: < Student />

\* Space (or) tab character is not allowed in the element name (or) in attribute name.

In the above simple XML Pgm "employees" is a root element. The root element is a element which is parent of all other elements in XML document.

The other elements are called contained elements. There can be one (or) more elements inside container elements. Such elements are called child elements. In the above example "employee" is a container element which contains the child elements "name" and "dept".

Attributes & values :-

Attributes are designed to contain date related to a specific element. Attributes must be enclosed within double quotes (or) single quote.

Ex: < Person gender = "female" >

XML Elements vs Attributes

Ex1: < Person gender = "female" >

< first name > Dutta < /first name >

< last name > Sai Eswari < /last name >

< / Person >

Ex2: < Person >

< gender > female < /gender >

< first name > Dutta < /first name >

</last name> Sai Eswari </last name>  
</person>.

In the "Ex1" gender is an attribute. In the "Ex2" gender is an element. Both examples provide same info.

### Name Spaces in XML

Sometimes we need to create two different elements by the same name. The XML document allows us to create different elements which are having the common name. This technique is known as namespace.

Ex: <file-Description>

<text fname = "input.txt">

<describe> It is text file </describe>

</text>

<text fname = "flower.jpg">

<describe> It is an image file </describe>

</text>

</file-Description>

The above document does not produce any error although the element text is used for two different attribute values.

### DOCUMENT TYPE DEFINITION (DTD)

The document type definition is used to define the basic building blocks of any XML document.

Using DTD we can specify various elements types, attributes and their relationship. Basically DTD is used to specify the set of rules for structuring

dates in any XML file.

If we want to put some information about students in XML file then, generally we use tag Student followed by his/her name, address, standard and marks. That means we are actually specifying the manner by which information should be arranged in XML file.

### Building Blocks of XML

1.) Elements :- Elements are the main building blocks of both XML & HTML. These are used for defining the tags. They typically consist of opening & closing tags. Mostly only one element is used to define a single tag.

Ex:- <body> Some text </body>  
<message> Some text </message>

2.) Attributes :- Attributes are generally used to specify the values of the element. They provide extra information about elements. These are specified with in double quotes.

Ex:- <flag type = "True">

The name of the element is "flag", name of the attribute is "type", value of the attribute is "True".

3.) Entities :- Some characters have special meaning in XML, like less than sign (<) that defines starting of XML tag.

Following entities are predefined in XML.

Entity Reference	Character.
&lt;	<
&gt;	>
&amp;	&
&quot;	"
&apos;	'

4.) CDATA: CDATA Stands for character data. CDATA is text that will not be parsed by a parser. Tags inside the text will not be treated as markup and entities will not be expanded.

5.) PCDATA: PCDATA means parsed character data. PCDATA is text that will be parsed by a parser. The text will be examined by parser for entities and markup. Tags inside the text will be treated as markup and entities will be expanded. Parsed character data should not contain any &, <, or > characters, these need to be represented by &amp; ; &lt; ; and &gt; ; entities.

- XML document contains data where as DTD files contains rules that apply to data.
- DTD must saved with '.dtd' extension.
- An XML file contains an reference to DTD file.
- A DOCTYPE declaration in an xml document specifies that we want to include a reference to DTD file.
- DTD defines doc structure with a list of legal elements & attributes.

## Syntax

```
<!DOCTYPE DOCUMENT [  
    <!ELEMENT ELEMENTNAME1 (attribute name)  
        appearance of attributes>  
    <!ELEMENT ELEMENTNAME2>  
    :  
    :  
    :  
    <!ELEMENT ELEMENTNAME n>  
<!AT LIST element-name attribute-name  
        attribute-type default-value>  
]
```

Types of DTD :- There are 2 types of DTD's.

- 1.) Internal DTD :- In this we enclose DTD in a `<DOC TYPES>` element providing the name of root element of document.

## Syntax:

```
<!DOCTYPE root-element [element-declaration]>
```

Root-element is the name of root element and element declaration is where you declare the elements.

Ex: <?xml version = "1.0"?>

```
<!DOCTYPE employee [
```

<!ELEMENT name (# PCDATA) >

<!ELEMENT dept (#PCDATA)>

75.

Employee

<name> Sai Eswari </name>

<dept> ST </dept>

<Employee>

External DTD :- In external DTD elements are declared outside the XML file. They are accessed by specifying the attributes which may be either the legal '.dtd' file (or) valid URL.

If the DTD is declared in an external file, the `<!DOCTYPE>` definition must contain a reference to DTD file.

Syntax:-

`<!DOCTYPE root-element SYSTEM "filename.dtd">`

Ex:- `<?xml version = "1.0"?>`

`<!DOCTYPE employee SYSTEM "emp.dtd">`

`<employee>`

`<name> Sai Eswari </name>`

`<dept> IT </dept>`

`</employee>`

`<!ELEMENT name (#PCDATA)>` } Save with

`<!ELEMENT dept (#PCDATA)>` } emp.dtd.

Uses of DTD

1.) with DTD each of your XML files can carry the description of its own format.

2.) with DTD independent groups of people can agree to use a standard DTD for interchanging data.

3.) your application can use standard DTD to verify that the data you received from outside world is valid.

4.) You can also use DTD to verify your own data.

5.) DTD's are relatively Simple and compact.

## Limitations of DTD

- They are not written in XML Syntax, which means you have to learn a new syntax in order to use.
- DTD's are not Object Oriented.
- DTD's don't support namespaces very well.
- we can't use multiple DTD to validate one XML doc.
- Large DTD's are hard to read and maintain.

## XML SCHEMAS

- The XML Schemas are used to represent the structure of XML document. The goal (or) purpose of XML Schema is to define the building blocks of an XML document. These can be used as an alternative to XML DTD. The XML schema language is called as XML Schema Definition (XSD) language. The XML Schema became the world wide web Consortium (W3C) recommendation in 2001.
- XML Schema defines elements, attributes, elements having child elements, order of child elements. It also defines fixed and default value of elements & attributes.
- XML Schema also allows the developer to use data types.

Ex:- Student schema.xsd.

```
<?xml Version = "1.0"?>
<xsd:schema xmlns:xsd = "http://www.w3.org/2001/XMLSchema">
    <xsd:element name = "student">
        <xsd:complexType>
            <xsd:sequence>
                <xsd:element name = "name" type = "xsd:string"/>
```

$\langle \text{xs:element name} = \text{"address"} \text{ type} = \text{"xs:string"} \rangle$   
 $\langle \text{xs:element name} = \text{"std"} \text{ type} = \text{"xs:string"} \rangle$   
 $\langle \text{xs:element name} = \text{"marks"} \text{ type} = \text{"xs:string"} \rangle$

$\langle \text{/xs:sequence} \rangle$

$\langle \text{/xs:complexType} \rangle$

$\langle \text{/xs:element} \rangle$

$\langle \text{/xs:schemas} \rangle$

Simple Type :- Simple type element is used only in the content of text. Some of predefined Simple types are :  
 $\text{xs:integer}$ ,  $\text{xs:boolean}$ ,  $\text{xs:string}$ ,  $\text{xs:date}$

Complex Type :- A complex type is a container for other element definitions. This allows you to specify which child elements an element can contain & to provide some structure within your XML documents.

$\langle \text{xs:element name} = \text{"Address"} \rangle$

$\langle \text{xs:complexType} \rangle$

$\langle \text{xs:sequence} \rangle$

$\langle \text{xs:element name} = \text{"name"} \text{ type} = \text{"xs:string"} \rangle$

$\langle \text{xs:element name} = \text{"Company"} \text{ type} = \text{"xs:string"} \rangle$

$\langle \text{xs:element name} = \text{"Phone"} \text{ type} = \text{"xs:int"} \rangle$

$\langle \text{/xs:sequence} \rangle$

$\langle \text{/xs:ComplexType} \rangle$

$\langle \text{/xs:element} \rangle$

Global Type :- With global type, you can define a single type in your document, which can be used by all other references. for example, suppose you want to generalize the Person & Company for different addresses of Company.

$\langle \text{xs:element name} = \text{"AddressType"} \rangle$

$\langle \text{xs:ComplexType} \rangle$

<xs: Sequence>  
 <xs: element name = "name" type = "xs:string"/>  
 <xs: element name = "Company" type = "xs:string"/>  
 </xs: Sequence>  
 </xs: Complex Type>  
 </xs: element>  
  
 <ns: element name = "Address 1">  
 <xs: Complex Type>  
 <ns: Sequence>  
 <xs: element name = "address" type = "Address Type"/>  
 <xs: element name = "phone1" type = "xs:int"/>  
 </xs: Sequence>  
 </ns: Complex Type>  
 </ns: element>  
 <ns: element name = "Address 2">  
 <xs: Complex Type>  
 <ns: Sequence>  
 <xs: element name = "Address" type = "Address Type"/>  
 <xs: element name = "Phone2" type = "xs:int"/>  
 </ns: Sequence>  
 </ns: Complex Type>  
 </xs: Sequence>  
 </ns: Complex Type>  
 </ns: element>

Advantages :-

- 1.) The Schemas are more specific.
- 2.) The Schema provide the support for datatypes.
- 3.) The Schema is aware of namespaces.

- 1) XML Schema allows the creation of complex and reusable Content models easily.

Advantages of Schema over DTD:-

- 1.) Both Schemas & DTD's are useful for defining structural Components of XML. But the DTD's are basic & cannot be much specific for complex operations.
- 2.) The Schemas provide support for defining type of data. The DTD's do not have this ability.
- 3.) The Schemas are namespace aware & DTD's are not.
- 4.) The XML Schema is written in XML itself & has a large number of built in & derived types.
- 5.) Large no. of web applications can be built using XML Schema. On the other hand relatively simple & compact operations can be built using DTD.

### DOCUMENT OBJECT MODEL (DOM)

DOM is a language-neutral & platform-independent object model used to represent XML documents. DOM helps scripts & programs to access, add, delete & edit content, structure and style of XML documents dynamically.

DOM is standardized by World Wide Web Consortium (W3C). The primary objective of this standard was to model HTML document in an object-oriented way so that it could be exposed to scripts, and scripts could access & manipulate HTML documents through this object model dynamically.

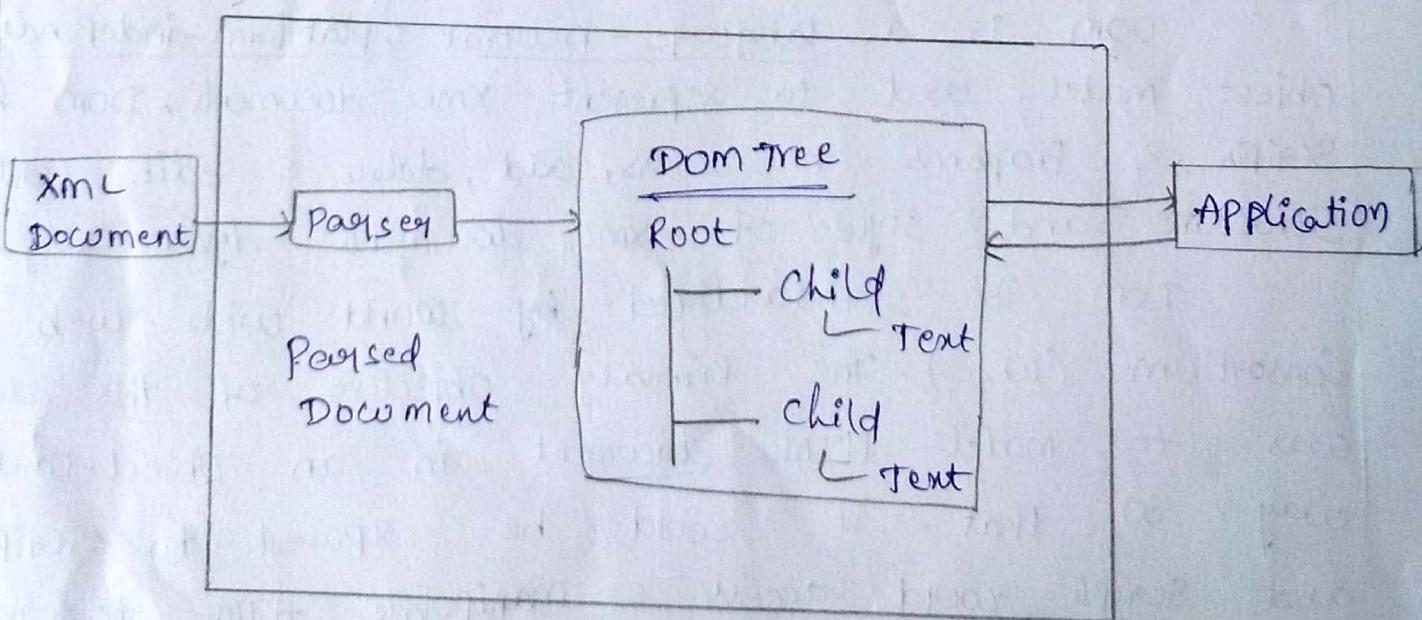
DOM defines the objects and Properties and methods (interface) to access all XML elements.

DOM specification is divided into 3 major parts:-

- 1.) DOM Core :- This portion defines the basic set of interfaces and objects for any structured document.
- 2.) XML DOM :- This part specifies the standard set of objects and interfaces for XML document only.
- 3.) HTML DOM :- This part defines the interfaces and objects for HTML documents only.

XML DOM is a standard object model for XML. XML documents have hierarchy of informational units called nodes. DOM is a standard programming interface of describing those nodes and the relationship between them.

As XML & DOM also provides an API that allows a developer to add, edit, move (or) remove nodes at any point on the tree in order to create an application.



Above is the diagram for the DOM structure, which depicts that Parser evaluates an XML document as DOM structure by traversing through each node.

A DOM document is a collection of nodes (or) pieces of information, organised in hierarchy. Some types of nodes may have child nodes of various types and others are leaf nodes that cannot have anything below them in the document. Below is a list of the nodes types, and which node types they may have as children.

- 1.) Document :- Only one Document type node exists for each XML document. It is essentially a container of all components of an XML document such as XML declaration, elements, attributes, comments, entities and so-on.
- 2.) Element :- This interface represents an elements in XML document. Elements may have attributes. It has several useful methods to get an Attr node (or) value of attribute by its name.
- 3.) Attr :- An Attr type node represents an attribute of an Element node ---- text, Entity Reference.
- 4.) Text :- It represents the textual content of Element type node. The text content of an element node is represented as separate Text type node.
- 5.) Document fragment :- The implementation of Document object can be heavyweight as a large no. of methods and properties have been defined for it---- element, Processing Instruction, comment, Text, CDATA Section, Entity Reference.
- 6.) Document Type :- It provides interfaces to get information about document, including list of entities defined for this document.

- 7.) Entity Reference :- A node of this type represents entity reference in the document. --- Element, Processing Instruction, Comment, Text, CDATA Section, Entity Reference
- 8.) Processing Instruction :- This interface represents a "processing instruction", which is used in XML to provide specific information about document - Processor.
- 9.) Comment :- It represents a comment in an XML document. Note that all the characters between starting '

10.) CDATA Section :- This type of node represents CDATA Section in XML doc. Note that no lexical check is done on content of CDATA Section.

11.) Entity :- It represents a known, either unparsed or parsed entity. --- Element, Processing Instruction, Comment, Text, CDATA Section, Entity Reference.

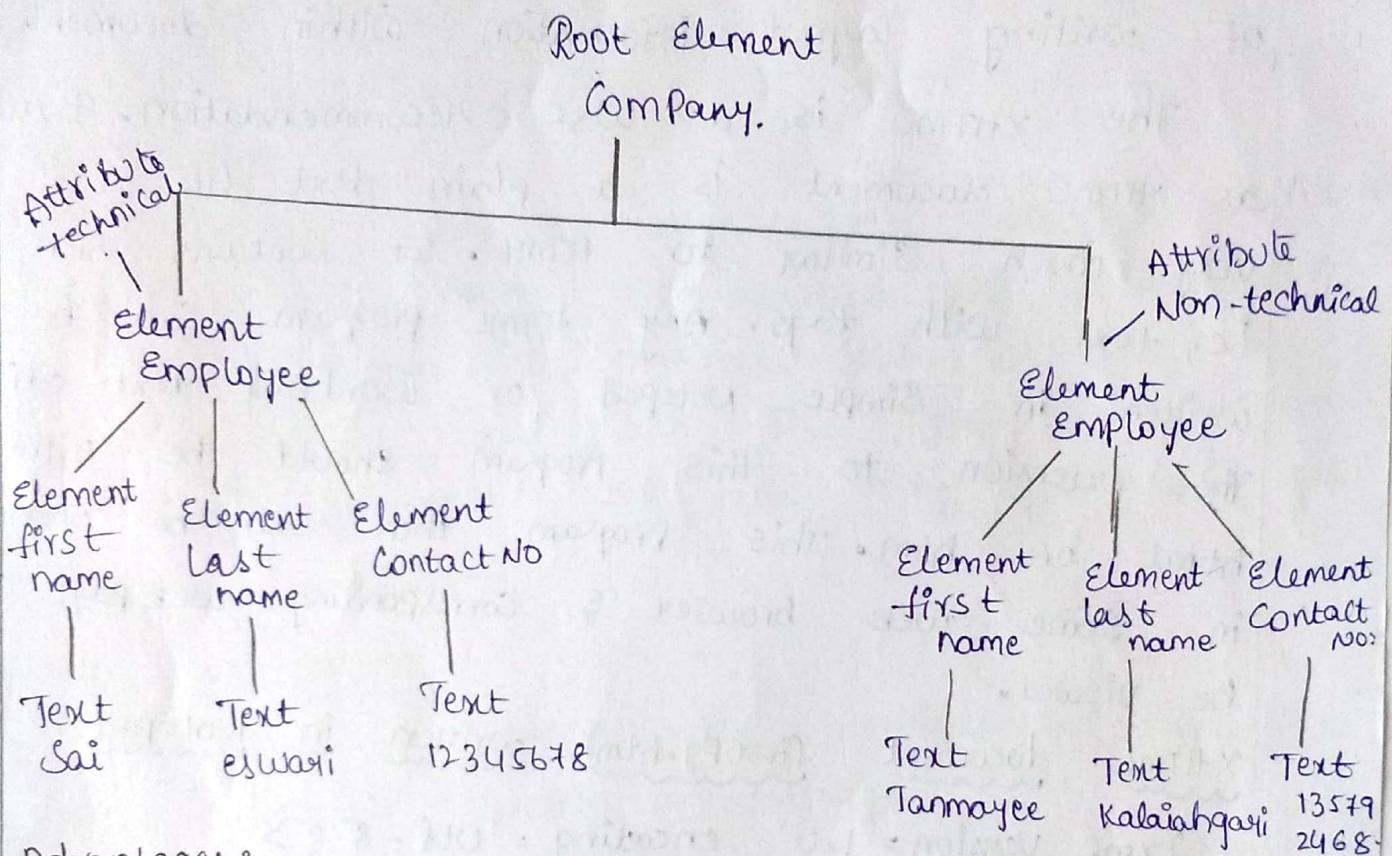
Eg:- Consider DOM representation of following XML doc node.xml.

```
<?xml version = "1.0"?>
<Company>
    <Employee category = "technical">
        <first name> Sai </first name>
        <last name> Eswari </last name>
        <Contact No:> 12345678 </Contact No:>
    </Employee>
    <Employee category = "non-technical">
        <first name> Tanmayee </first name>
        <last name> Kalaiyahgari </last name>
        <Contact No> 135792468 </Contact No>
```

< / Employee >

< / Company >

The Document Object Model of above XML document would be as follows:



Advantages :-

- XML DOM is language and Platform independent.
- XML DOM is traversable - Information in XML DOM is organised in hierarchy which allows developer to navigate around hierarchy looking for specific information.
- XML DOM is Modifiable - It is dynamic in nature providing developer a scope to add, edit, move (or) remove nodes at any point on tree.

Disadvantages :-

- It consumes more memory.
- Due to larger usage of memory its operational speed, compared to SAX is slower.

## XHTML

XHTML stands for extensible Hypertext markup language. Hypertext is simply a piece of text which works as a link. Markup language is language of writing layout information within documents.

The XHTML is a W3C recommendation. Basically an XHTML document is a plain text file & it is very much similar to HTML. It contains rich text i.e. text with tags. Any HTML program can be written in simple Notepad or Word Pad text editors. The extension to this program should be either .html or .htm. This program then can be opened in some web browser & corresponding web page can be viewed.

XHTML document firstPg.html written in Notepad :-

```
<?xml version="1.0" encoding="Utf-8"?>
```

```
<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional //EN".
```

```
"http://www.w3.org/TR/xhtml/DTD/xhtml1.dtd">
```

```
<head>
```

```
    <title> My Page </title>
```

```
<body>
```

```
    <h2> This is my first web page </h2>
```

```
</body>
```

```
</html>
```

Every XHTML document must begin with an XML declaration element that simply identifies the document because it is based on XML. Thus

In this line the XML version is specified as 1.0 & encoding as utf i.e., unicode encoding.

Note that XHTML document consists of some strings enclosed within angular brackets. Such strings are called tags. Hence various tags that are written in above doc are `<html>`, `<h2>`, `<head>` & `<body>`.

Generally code for any scripting language is written in some text editor like Notepad. While saving XHTML document written in Notepad, save it using file extension `'.html'` and for file type select option `'All files'` otherwise your file may get saved as `'first Pg. html.txt'`.

### HTML vs XHTML

XHTML is syntactically identical to HTML. But XHTML follows certain restrictions.

#### HTML:-

- 1.) The HTML tags are case insensitive. Hence `<body>` or `<BODY>` or `<Body>` are treated as one & the same.
- 2.) We can omit the closing tags sometimes in HTML document.
- 3.) In HTML the attribute values is not always necessary to quote. In fact numeric attribute values are rarely quoted in HTML. Only when special characters (or) white spaces are present in attribute values then only it is essential to put quotes around them in HTML.

- 4.) In HTML there are some implicit attribute values and standard attributes.
- 5.) In HTML even if we do not follow the nesting rules strictly it does not cause much difference.

### XHTML :-

- 1.) The XHTML is case sensitive and all the tags in XHTML document must be written in lowercase.
- 2.) For every tag there must be closing tag. Some browsers get confused if the closing tag is not given. There are two ways by which we can mention the closing tags.

`<a href = "firstPage.html"> </a>`  
(Or)

`<a href = "firstPage.html"/>`

- 3.) In every XHTML document the attribute values must be quoted.
- 4.) In every XHTML the attribute values must be specified explicitly.
- 5.) In XHTML document the nesting rules must be followed. These nesting rules are:
- A form element cannot contain another form element.
  - An anchor element does not contain another form element.
  - List elements cannot be nested in the list element.
  - If there are two nested elements then the inner element must be enclosed first before closing outer element.
  - Text elements cannot be directly nested in form elements.

## Standard Structure :-

An XHTML document consists of 3 main parts.

1) DOCTYPE declaration.

2) <head> section.

3) <body> section.

```
<!DOCTYPE >  
<html>  
  <head>  
    <title>----</title>  
  </head>  
  <body>  
  ----  
</body>  
</html>
```

The DOCTYPE declaration should be

<!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional //EN"  
"http://www.w3.org/TR/xhtml1/DTD/xhtml11.dtd">.

DOCTYPE specifies the document type. The Document types is specified by DTD. The XHTML Syntax rules are specified by file xhtml 11.dtd file.

There are 3 types of XHTML DTD's and those along with their uses are given below.

1.) XHTML 1.0 Strict :- when we want a clean markup code then this type of dtd is used.

2.) XHTML 1.0 Transitional :- when we want to use some HTML features in the existing XHTML document.

3.) XHTML 1.0 frameset:- When we want to make use of frames in the XHTML document.

PARSING XML DATA: DOM and SAX PARSERS.

The primary goal of an XML processor is to parse the given XML document. Java has a rich source of in-built API's for Parsing the given XML document. It is parsed in two ways.

- Tree based parsing.
- Event based parsing.

As we know that DOM is used to parse the XML document using the tree structure. We can access the information of XML document by interacting with the tree nodes. On the other hand Simple API for XML (i.e., SAX) allows us to access the information of XML document using sequence of events. Thus there are two methods of parsing the XML document.

- Parsing using DOM (tree based).
- Parsing using SAX (event based).

Difference between DOM and SAX.

SAX

DOM

- |   |   |
|---|---|
| 1.) Simple API for XML                    | 1.) Document object models.                                     |
| 2.) It reads an XML file as node by node. | 2.) It stores entire XML doc.                                   |
| 3.) It detects an events                  | 3.) It builds an in-memory tree representation on XML document. |

- w) It informs the event to our application program.
- x) In SAX, application Pgm acts upon the event.
- y) Parsing continues till all the events (or) end of XML is parsed.
- z) It uses following package  
import javax.xml.parsers;
- a) It parses the whole XML document at a time.
- b) It uses more amount of memory.
- c) It uses following package  
import java.xml.parsers.\*;

DOM & SAX Pgm example.

### Dom.java

```

import java.io.*;
import java.xml.parsers.*;
import org.w3c.dom.*;
import org.xml.sax.*;

public class DOM
{
    static public void main (String [] args)
    {
        try {
            System.out.print("enter name of document");
            BufferedReader input = new BufferedReader (new InputStreamReader (System.in));
            String file_name = input.readLine();
            file fp = new file (file_name);
        }
    }
}

```

```
if (fp.exists())  
{  
    try {  
        DocumentBuilderFactory factory_obj = Document  
       BuilderFactory.newInstance();  
        DocumentBuilder builder = factory_obj.newDocument  
        Builder();  
  
        InputSource ip_src = new InputSource (file_name);  
        Document doc = builder.parse (ip_src);  
        S.O. PIn (file_name + " is well-formed");  
    }  
    catch (Exception e) {  
        S.O.Pln (file_name + " is not well-formed");  
        System.exit (1);  
    }  
}  
else {  
    S.O.P ("file not found");  
}  
}  
catch (IOException ex)  
{  
    ex.printStackTrace();  
}  
}
```

Student.xml

```

<?xml version = "1.0"?>
<student>
    <Roll No> 10 </Roll No>
    <Personal info>
        <Name> Sai </Name>
        <Phone> 123456 </Phone>
    </Personal info>
    <Class> 10th </Class>
    <Sub> Maths </Sub>
    <Marks> 98 </Marks>
    <Roll no> 20 </Roll no.>
    <Personal info>
        <Name> Parvathi </Name>
        <Phone> 789012 </Phone>
    </Personal info>
    <Class> 10th </Class>
    <Sub> Eng </Sub>
    <Marks> 100 </Marks>
</Student>

```

Output:

D:\SAX-example>javac Dom.java.

D:\SAX-example>java DOM.

Enter name of document: student.xml.

Student.xml is well-formed!

D:\SAX-example>