

Tema 1. Introducción a los lenguajes de marcas. XML

1. Lenguajes de marcas

Un **lenguaje de marcas** es un sistema de representación de información basado en etiquetas que permite estructurar, identificar y describir datos. La idea central es separar el contenido de su presentación, para favorecer la reutilización, el intercambio y el procesamiento automático.

Los aspectos clave que distinguen a los lenguajes de marcas son:

- Uso de etiquetas y estructura jerárquica (en forma de árbol).
- Separación entre contenido y presentación.
- Legibilidad por humanos y máquinas.
- Independencia de plataforma y formato (texto plano).



1.1. Clasificación de los lenguajes de marcas

Los lenguajes de marcas se pueden clasificar según su propósito o finalidad en:

- Presentación: HTML, Markdown, RTF, MathML.
- Intercambio de datos: XML, JSON, YAML.
- Específicos de dominio: SVG, MathML, LaTeX.

Ejemplo de un mismo dato representado en diferentes lenguajes:

HTML: `<p>Hola mundo</p>`

XML: `<mensaje>Hola mundo</mensaje>`

JSON: `{ "mensaje": "Hola mundo" }`

2. XML



XML proviene de **GML** (Generalized Markup Language), un lenguaje que surgió en la década de los 70 en el seno de la empresa IBM, con el objetivo de gestionar grandes cantidades de información. Dada su gran utilidad, en 1986, la Organización Internacional de Estándares lo normalizó creando **SGML** (Standard Generalized Markup Language).

Tras el nacimiento de la web, la popularización de HTML hizo este lenguaje de marcas creciera sin orden alguno, hasta que en 1996 el W3C se centró en crear un estándar que deberían adoptar todos los navegadores, para así facilitar el trabajo de los desarrolladores. XML nació como parte del estándar SGML, con el objetivo de solventar las carencias de HTML en lo que a tratamiento de la información se refiere. Esto es, cuando trabajamos con HTML, no sólo el contenido se mezcla con los estilos, dependiendo la presentación del navegador que se utilice, sino que este estándar no nos permite el intercambio de información con otros dispositivos.

XML (eXtensible Markup Language) se define como un metalenguaje extensible de etiquetas desarrollado por el W3C, como una simplificación y adaptación de SGML, permitiendo a los desarrolladores definir la gramática de lenguajes específicos. Propuesto como un estándar para el intercambio de información estructurada entre diferentes plataformas, XML no es realmente un lenguaje en particular, sino una manera de definir lenguajes para diferentes necesidades. Su papel en la actualidad es bastante importante, ya que permite la compatibilidad entre sistemas para compartir la información de una manera segura, fiable y fácil, además de utilizarse como base para la creación de lenguajes estándares que describen diferentes tipos de datos: **SVG** (*Scalable Vector Graphics*), **MathML** (*Mathematical Markup Language* – descripción de fórmulas matemáticas), **ODT** (*Open Document Format*), **RSS** (*Really Simple Syndication*), **ePUB** (formato de libro digital), **XHTML** (versión extendida de HTML que estudiaremos próximamente), etc.

3. Documentos XML

XML es un lenguaje de marcas orientado a la descripción, esto es, determinadas marcas permitirán dar significado al texto sin indicar cómo representarlo en pantalla. De esta forma, un documento XML contiene datos que se autodefinen, en tanto que en un documento HTML, como veremos más adelante, datos y representación están indisolublemente unidos. XML busca dar solución al problema de expresar información estructurada, de la manera más abstracta y reutilizable posible. Que la información sea estructurada quiere decir que se compone de partes bien definidas, y que esas partes se componen a su vez de otras partes. De esta manera, XML se puede emplear como un lenguaje limpio que permite describir información de forma precisa, sin tener que preocuparnos en ningún momento de su representación final. Veamos un sencillo ejemplo de documento XML.

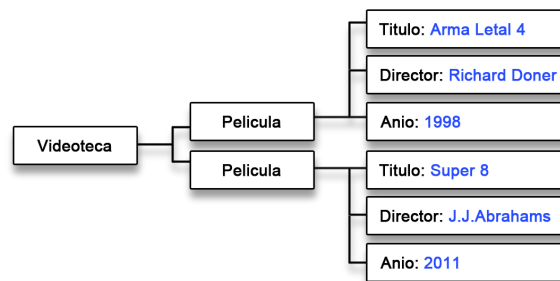
```
<?xml version="1.0" encoding="UTF-8"?>
<!DOCTYPE videoteca SYSTEM "videoteca.dtd">
<videoteca>
  <pelicula>
    <titulo>Arma Letal 4</titulo>
    <director>Richard Donner</director>
```

```

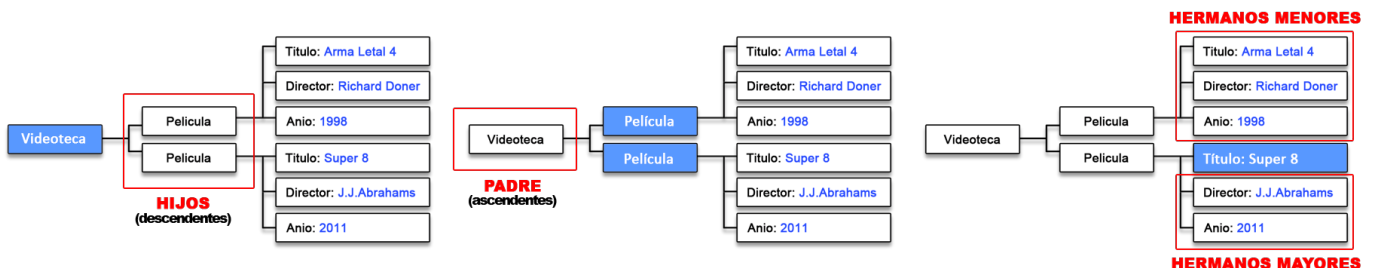
    <anio>1998</anio>
  </pelicula>
<pelicula>
  <titulo>Super 8</titulo>
  <director>J.J.Abrahams</director>
  <anio>2011</anio>
</pelicula>
</videoteca>

```

En definitiva, un documento XML está constituido por texto plano y **etiquetas** (marcas) que permiten clasificar al texto dotándolo de cierto significado. A diferencia de lo que ocurre con HTML, aquí no existen etiquetas predefinidas, por lo que el desarrollador podrá definir aquellas que desee según sus necesidades. Como vemos en el ejemplo anterior, la estructura del documento nos permite agrupar bajo diferentes etiquetas, creadas por el propio desarrollador, un conjunto de datos relacionados entre sí, generando una **estructura de tipo árbol** como la que vemos en la siguiente imagen.



La jerarquía formada entre los elementos del árbol establece automáticamente relaciones entre padres, hijos, hermanos, ascendentes y descendientes. De esta forma, las partes del árbol que tienen hijos se las denomina **nodos intermedios**, en tanto que aquellas que no tienen se conocen como **hojas**.



Es importante conocer la sintaxis de XML para poder crear documentos bien-formados. Un documento escrito bajo estándar XML debe seguir una estructura estrictamente jerárquica, en lo referente a las etiquetas que delimitan a sus elementos.

```
<p>El lenguaje HTML <b>permite <i>esto</b></i></p>
```

Según el estándar XML, la línea de código anterior sería incorrecta ya que éste exige que cada etiqueta esté correctamente contenida dentro de otra, teniendo en cuenta que todas deberán estar cerrarse en orden inverso al que se abrieron.

```
<p>En el lenguaje XML <b>la estructura <i>debe ser</i> jerárquica</b></p>
```

2.1. Elementos básicos y atributos

Como se introdujo al principio de la página, los **elementos básicos** de XML son las marcas o etiquetas. Éstas no son más que contenedores de información y deben aparecer en parejas, delimitando de esta manera su

contenido que, a su vez puede ser vacío, contener texto, atributos u otros elementos de XML. Las etiquetas pueden completarse con **atributos** que nos permitan definir ciertas características o propiedades de dicho elemento. Éstos están constituidos por pares *nombre-valor* que son definidos en la etiqueta de apertura.

etiqueta de apertura
etiqueta de cierre
`<distancia unidades="km"> 70 </distancia>`
atributo y su valor
contenido

Por ejemplo, supongamos que en el ejemplo propuesto anteriormente, deseamos indicar la calidad de cada una de las películas. Podemos hacerlo incorporando un atributo llamado `calidad`, al elemento `película`. De esta manera, el código XML anterior quedaría.

```
<pelicula calidad="Muy buena">
  <titulo>Arma Letal 4</titulo>
  <director>Richard Donner</director>
  <anio>1998</anio>
</pelicula>
```

No hay una regla fija para saber cuándo añadir atributos pero, algo que nos puede ayudar es pensar qué datos son reales y cuáles describen a los primeros. Los siguientes ejemplos contienen la misma información:

```
<carta fecha="2008-01-10">
  <a>Carlos</a>
  <de>Alba</de>
</carta>
```

```
<carta>
  <fecha>2008-01-10</fecha>
  <a>Carlos</a>
  <de>Alba</de>
</carta>
```

```
<carta>
  <fecha>
    <anio>2008</anio>
    <mes>01</mes>
    <dia>10</dia>
  </fecha>
  <a>Carlos</a>
  <de>Alba</de>
</carta>
```

Algunas cosas a tener en cuenta al utilizar atributos son:

- Los atributos **no** pueden contener múltiples valores (los elementos sí)
- Los atributos **no** pueden contener estructuras de árbol (los elementos sí)
- Los atributos **no** son fácilmente expandibles (para cambios futuros)

En XML también podemos encontrar **elementos vacíos**, estos son los que no tienen contenido y que, aunque pueden tener atributos, se abren y cierran con una sola etiqueta.

`<separador distancia="7" />`

```
<separador distancia="7"></separador>
```

Los elementos se emplean por tanto para representar jerarquías o contenido y, el orden en el que aparecen no es representativo, pudiendo encontrar a lo largo del documento XML múltiples ocurrencias de dicho elemento. Además, un elemento, etiqueta o marca puede contener atributos que modifican la información, no pudiendo aparecer un mismo atributo más de una vez en una misma etiqueta.

3.2. Cabecera del documento

Aunque no es obligatorio, es aconsejable que un documento XML comience con una línea que describa la versión de la especificación utilizada y el tipo de documento. En ocasiones también se incluye en esta definición el atributo **standalone**, cuyo valor puede ser **yes** o **no**, y especificará si un documento depende o no de otros, como una DTD.

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

Seguidamente, hemos de realizar la declaración que define qué tipo de documento estamos creando, para que pueda ser procesado correctamente. Es decir, indicamos una **DTD (Document Type Definition)** válida, donde se definen las reglas que debe cumplir nuestro documento XML. En nuestro caso, si nos fijamos en el ejemplo anterior, indicamos que el contenido de nuestro archivo XML se basa en la estructura definida en la DTD contenida en el archivo `videoteca.dtd`, cuyo elemento raíz es `videoteca`.

```
<!DOCTYPE videoteca SYSTEM "videoteca.dtd" >
```

A continuación, en el documento encontramos una estructura jerárquica de elementos que, a su vez pueden contener a otros elementos, estar asociados directamente con cierto contenido, o bien ser elementos vacíos. Un elemento de contenido es, por ejemplo:

```
<pelicula>Arma letal 4</película>
```

XML es un lenguaje estricto y es sensible a mayúsculas y minúsculas, por lo que si se define un elemento de una u otra manera, siempre tendremos que referirnos a él de la misma forma. Al emplear XML será necesario asignar nombre a las estructuras, tipos de elemento, entidades, etc. En este caso, debemos seguir lo exigido por el estándar que nos dice que: **“Un nombre debe empezar con una letra, o un guión bajo, continuándose con letras, dígitos, guiones, rayas, dos puntos o puntos.”** Está terminante prohibido que el nombre de un elemento empiece por la cadena XML en alguna de sus variantes (empleando mayúsculas y/o minúsculas). Por lo tanto, las principales características de XML, que deberemos tener muy en cuenta a la hora de crear documentos bien formados son:

- a. Todos los elementos XML deben cerrarse.

```
<p>This is a paragraph.</p>  
<br />
```
- b. XML es sensible a mayúsculas y minúsculas.

```
<mensaje>This is correct</mensaje>  
<mensaje>This is NOT correct</Mensaje>
```
- c. Los elementos XML deben estar siempre correctamente anidados.

```
<b><i>This text is bold and italic</i></b>
```
- d. Un documento XML tiene **un único elemento raíz**.

```
<raiz>
```

```
<hijo>
  <subhijo>.....</subhijo>
</hijo>
</raiz>
```

- e. Todos los atributos de un elemento XML deben encerrar su valor entre comillas.

```
<carta date="12/11/2007">
  <a>Carlos</a>
  <de>Alba</de>
</carta>
```

- f. Los comentarios en XML comienzan por `<!--` y terminan por `-->`. No pueden aparecer antes de la especificación del documento, versión y tipo de codificación.

```
<!-- This is a comment -->
```

- g. Los espacios en blanco (tabulador, nueva línea, retorno de carro y espacio) se conservan en las secciones **PCDATA** (*Parsed Character Data*) que encontramos como contenido de las etiquetas. Como valor de un atributo, los espacios en blanco adyacentes se condensarán en uno sólo.

XML:	Hello	Tove
HTML:	Hello Tove	

- h. XML utiliza 5 entidades predefinidas para representar los caracteres: ' (*'*), " (*"*), < (*<*), > (*>*), & (*&*).

```
<message>salary &lt; 1000</message>
```

3.3. Parser

Se conoce con el nombre de **parser** o **analizador XML**, a la aplicación que se encarga de leer un documento XML y determinar la estructura y propiedades de los datos contenidos en él, generando el árbol jerárquico asociado. Definimos un **parser validador** como aquel que, además se encarga de comprobar que el código XML que conforma el documento, cumple unas determinadas reglas propuestas por el desarrollador, comprobando de esta forma la semántica del documento e informando de posibles errores. Son muchas las herramientas existentes para el manejo de XML y sus tecnologías (DTD, XSD, XPATH, XSLT, etc.), entre todas ellas destacamos: XMLSpy, <oXygen/>, XML Copy Editor, XMLPad Pro Edition, Exchanger XML Editor y Liquid XML Editor.

3.4. Espacio de nombres

En ocasiones, en un mismo documento XML, se producen conflictos de nombres entre elementos o atributos con diferentes definiciones. Esto supone un problema para el parser que no sabría cómo manejar etiquetas iguales. Los **espacios de nombres** o **namespacing** evitan este problema, indicando en cada etiqueta el contexto para cada una de ellas, permitiendo de esta manera el poder diferenciarlas.

```
<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE videoteca SYSTEM "videoteca.dtd">
<videoteca>
```

```

<titulo>Videoteca de Alex</titulo>
<pelicula>
  <titulo>Arma Letal 4</titulo>
  <director>Richard Donner</director>
  <anio>1998</anio>
</pelicula>
</videoteca>

```

Observamos en el ejemplo anterior que la etiqueta **titulo** se utiliza en contextos diferentes: el primero de ellos da título a la videoteca y el segundo se corresponde con el título de la película en particular. Esto, como se dijo anteriormente, supone un conflicto para el parser. La solución a este problema consiste en anteponer a la etiqueta un nombre (espacio de nombres) que identifique el propietario de la misma.

```

<?xml version="1.0" encoding="UTF-8" ?>
<!DOCTYPE videoteca SYSTEM "videoteca.dtd">
<videoteca>
  <video:titulo>Videoteca de Alex</video:titulo>
  <video:pelicula>
    <ficha:titulo>Arma Letal 4</ficha:titulo>
    <ficha:director>Richard Donner</ficha:director>
    <ficha:anio>1998</ficha:anio>
  </video:pelicula>
</videoteca>

```

En resumen, un espacio de nombres supone una recomendación del W3C para eliminar ambigüedades entre elementos y atributos de un documento XML. Generalmente, estos espacios de nombres suelen tener como identificador único una **URI** (identificador único de recurso), que define un nombre lógico para el espacio de nombres, resolviendo de esta manera posibles problemas de duplicidades. Utilizamos el atributo **xmlns** (xml namespaces) que permitirá asignar un espacio de nombres a un prefijo en el documento.

```

<raiz xmlns:h="http://www.w3.org/TR/html4/"
xmlns:f="https://www.w3schools.com/furniture">

```

```

<h:table>
  <h:tr>
    <h:td>Apples</h:td>
    <h:td>Bananas</h:td>
  </h:tr>
</h:table>

<f:table>
  <f:name>African Coffee Table</f:name>
  <f:width>80</f:width>
  <f:length>120</f:length>
</f:table>

</raiz>

```

El **ámbito** de declaración de un espacio de nombres cubre el elemento donde se ha declarado y todos sus descendientes.

```

<videoteca xmlns:video="videoteca:video">

```

```

        xmlns:ficha="videoteca:ficha">
<video:titulo>Videoteca de Alex</video:titulo>
<ficha:película>
    <ficha:titulo>Arma Letal 4</ficha:titulo>
    <ficha:director>Richar Donner</ficha:director>
    <ficha:anio>1998</ficha:anio>
</ficha:película>
</video:videoteca>

```

Los atributos pueden pertenecer, o no, al mismo espacio de nombres al que pertenece el elemento al que está asociado. Tal y como podemos ver en el siguiente ejemplo, el atributo `trabajo:empno` del elemento `<emp:datos>` pertenece a un espacio de nombres diferente a éste. Sin embargo, el atributo `deptno` no pertenece a ningún espacio de nombres.

```

<empleado xmlns:emp="empresa:emp">
    <emp:datos trabajo:empno="1234"
        xmlns:trabajo="empresa:trabajo">
        <emp:departamento deptno="10">Ventas</emp:departamento>
        <emp:nombre>Roberto</emp:nombre>
        <emp:apellidos>Ramírez</emp:apellidos>
    </emp:datos>
</emp:empleado>

```

Lógicamente, si dos atributos pertenecen a espacios de nombres diferentes, aunque tengan el mismo nombre se considerarán distintos, por lo que podrán asociarse con el mismo elemento.

```

<emp:datos trabajo:empno="1234" emp:empno="e_1234"> ... </emp:datos>

```

En caso de que las etiquetas de un documento pertenezcan mayoritariamente a un mismo espacio de nombres, lo lógico es indicar cuál es el namespace por defecto. En este caso, el ámbito de aplicación del espacio de nombres es del elemento en el que se ha declarado y sus elementos descendientes, pero **no sus atributos**. Para el ejemplo de la videoteca, supongamos que el namespace por defecto es `video`, así que eliminamos este sufijo del atributo `xmlns` en su definición, tal y como vemos a continuación:

```

<videoteca xmlns="videoteca:video"
    xmlns:ficha="videoteca:ficha">

```

4. Otros lenguajes XML

Al comienzo del tema mencionamos algunos lenguajes, basados en XML, que se utilizan para propósitos específicos. Entre ellos encontramos los siguientes:

- a. **SVG (Scalable Vector Graphics)**. Auspiciado por el W3C, SVG se define como un lenguaje de marcas para la representación de gráficos vectoriales bidimensionales, que actualmente están soportados por cualquier navegador. Almacena gráficos vectorizados escalables en formato texto, lo que, no sólo permite editarlo con cualquier aplicación de edición de textos, sino que también da lugar a archivos muy compactos. Además, soporta hojas de estilo y pueden generarse dinámicamente en un servidor web, como respuesta a la interacción con el usuario. Este formato es utilizado por numerosas aplicaciones gráficas, como Adobe Illustrator, Corel Draw, Inkscape, Mayura Draw, etc.

- b. WML (Wireless Markup Language).** Es un lenguaje de marcas utilizado para representar la información que se visualiza en dispositivos móviles y asistentes digitales que utilicen tecnología **WAP** (Wireless Application Protocol - estándar utilizado para comunicaciones inalámbricas).

```
<?xml version="1.0" ?>
<!DOCTYPE wml PUBLIC "-//WAPFORUM//DTD WML 1.1//EN"
"http://www.wapforum.org/DTD/wml_1.1.xml">

<wml>
  <card id="MiPrimerWML" title="Primer WML">
    <p align="center">
      Mi primer ejercicio en WML
    </p>
  </card>
</wml>
```



- c. RSS (Really Simple Syndication).** Se utiliza para syndicar o compartir contenido en la web, difundiendo información actualizada a los usuarios que se hayan suscrito previamente a la fuente de contenidos. Un archivo RSS contiene información básica sobre las novedades del sitio, como título, fecha de publicación o descripción. Generalmente se utilizan agregadores, esto es, aplicaciones web o de escritorio, diseñadas específicamente para la lectura de dichos contenidos, encargadas además de aplicar un estilo al contenido, presentándolos de forma atractiva al usuario. Alguno de los agregadores más utilizados son Flipboard, Feedly, Menéame, Scoop.it, Divúlgame, etc.
- d. DocBook.** Además de una aplicación, es un lenguaje utilizado para crear documentación, generalmente de tipo técnico, utilizando un formato neutro independiente de su presentación que permite especificar tanto contenido, como la estructura.

```
1 <?xml version="1.0" encoding="UTF-8" standalone="no"?>
2 <!DOCTYPE bookarticlearticlearticlearticle PUBLIC "-//OASIS//DTD DocBook XML V4.5//EN" "../
3 docbook-xml-4.5/docbookx.dtd">
4 <book>
5   == Ejemplo Docbook Book
6   <chapter>
7     == Primer capítulo
8     == Primera seccion del primer capitulo
9     Ipsum reversus ab viral inferno, nam rick grimes malum cerebro. De carne lumbering
10    animata corpora quaeritis. Summus brains sit, morbo vel maleficia? De apocalypsi
11    gorger omero undead survivor dictum mauris...
12    == Segunda seccion del primer capitulo
13    Hi mindless mortuis soulless creaturas, imo evil stalking monstra adventus resi
14    dentevil vultus comedat cerebella viventium. Qui animated corpse, cricket bat max
15    brucks terribilem incessu zombie...
16  </chapter>
17 </book>
```

Bibliografía

- Lenguajes de Marcas. XML. Validación de documentos
Jorge Sánchez Asenjo
- Lenguajes de Marcas y Sistemas de Gestión de la Información
Juan Manuel Castro Ramos, José Ramón Rodríguez Sánchez
Editorial Garceta
- Internet Encyclopedia
Hossein Bidgoli
California State University
- W3Schools