

Teoría de Lenguajes

Primer cuatrimestre 2013

Departamento de Computación
Facultad de Ciencias Exactas y Naturales
Universidad de Buenos Aires

TP Micro HTML Prettyprint

Integrante	LU	Correo electrónico
Ezequiel Gutesman	715/02	egutesman@gmail.com
Mauricio Alfonso	65/09	mauricioalfonso88@gmail.com
Victor Hugo Montero	707/98	vico.walker@gmail.com

1 Introducción

Cuando comenzamos a pensar en cómo encarar el problema del desarrollo de la gramática, pensamos que podríamos tratar cada carácter (letra, número, espacio, etc.) como un símbolo terminal diferente, pero terminaríamos con una cantidad inmanejable de símbolos no terminales y producciones, es decir con una gramática enorme. Por esa razón definimos una serie de Token Léxicos que al pasárselos a un analizador léxico transforman el texto de entrada de una serie de caracteres en una serie de tokens.

Definimos los tokens usando expresiones regulares. Primero definimos un token por cada tag de HTML y por último un token para todo el texto que sobre. El analizador léxico tokeniza el archivo buscando los tokens en orden, de manera que todo el texto que halla entre dos tags queda tokenizado como `texto`.

El caso de la sección `script` de HTML nos presentó un problema: entre los tags `<script>` y `</script>` puede haber cualquier combinación de caracteres excepto el tag de cierre de script. Esto quiere decir que podría haber incluso otros tags HTML dentro de un script, que deberían ser ignorados y que podrían no formar HTML válido. Pensamos en un principio en tokenizar el texto interior de los scripts con el analizador léxico como `texto_sin_script` para tomar todo el texto que no tenga el tag `</script>` como un token. Sin embargo no es posible tokenizar el texto de esta manera, ya que no tiene sentido poner el token `texto_sin_script` antes ni después de los tags de HTML. Esto se debe a que si el analizador léxico busca primero el token `texto_sin_script` y luego los tags, el primer token estaría “absorbiendo” todos los tags excepto `</script>` como parte de un script, incluso cuando no forman parte de un script; entonces no tokenizaría nunca los tags HTML. Si en cambio tokenizáramos primero los tags HTML y después `texto_sin_script`, los tags del interior de los scripts serían reconocidos y el tag `texto_sin_script` no cumpliría su función. Por esta razón decidimos tokenizar únicamente `texto` para todo el texto que sobre luego de haber reconocido los tags HTML, y reconocer los tags que pueda haber en el script desde la gramática y no desde el analizador léxico.

Además de tokenizar, el analizador léxico también se encarga de eliminar los comentarios y los espacios en blanco consecutivos, ya que estos podrían estar en cualquier parte del archivo recibido y debemos ignorarlos por completo.

En la gramática definimos un símbolo no terminal S para todo el documento.

El símbolo no terminal H contiene el interior del HTML, que a su vez puede tener un $HEAD$ y un $BODY$.

Los símbolos $HEAD$ y $BODY$ representan las secciones head y body de HTML respectivamente, y los símbolos HE y B representan el interior de dichas secciones.

Con el símbolo *SCS* definimos una serie de cero o más scripts, que identificamos con el símbolo *SC*.

El símbolo *TSC* representa el texto del interior de un script. Este es el símbolo con más producciones, ya que puede tener en su interior cualquier tag excepto `</script>` y además estos pueden ir de cualquier manera, sin formar HTML válido.

El interior de la sección body lo definimos recursivamente con el símbolo *B*, entendiendo que entre cada par de tags de apertura y cierre del mismo tipo puede haber más texto HTML válido.

Decidimos que sólo puede haber un elemento `<TITLE>...</TITLE>` dentro del head.

2 Tokens Léxicos

Los tokens léxicos están descritos por expresiones regulares. La siguiente tabla describe para cada token su expresión regular correspondiente.

<i>Token Léxico</i>	<i>Expresión Regular</i>
<HTML>	<html>
</HTML>	</html>
<HEAD>	<head>
</HEAD>	</head>
<BODY>	<body>
</BODY>	</body>
<TITLE>	<title>
</TITLE>	</title>
<SCRIPT>	<script>
</SCRIPT>	</script>
<DIV>	<div>
</DIV>	</div>
<P>	<p>
</P>	</p>
texto	.*

Antes de tokenizar, el analizador léxico elimina espacios consecutivos y comentarios. Los comentarios cumplen con la expresión regular <!--.*-->

3 Gramática

La siguiente gramática define el MicroHTML:

$$G = \langle N, \Sigma, P, S \rangle$$

Donde:

$$\begin{aligned} N &= \{S, H, HEAD, BODY, HE, SCS, SC, TSC, B\} \\ \Sigma &= \{texto, \langle HTML \rangle, \langle /HTML \rangle, \langle HEAD \rangle, \langle /HEAD \rangle, \\ &\quad \langle TITLE \rangle, \langle /TITLE \rangle, \langle SCRIPT \rangle, \langle /SCRIPT \rangle, \\ &\quad \langle BODY \rangle, \langle /BODY \rangle, \langle H1 \rangle, \langle /H1 \rangle, \langle DIV \rangle, \langle /DIV \rangle, \\ &\quad \langle P \rangle, \langle /P \rangle, \langle BR \rangle\} \\ S &= S \end{aligned}$$

Que de acuerdo a la clasificación de Chomsky es una gramática *tipo 0*, o *sin restricciones*. Las producciones (P) se detallan en la siguiente subsección.

3.1 Producciones - P

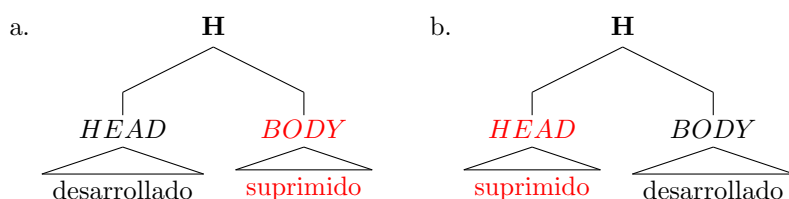
$$\begin{aligned} S &\rightarrow \langle HTML \rangle H \langle /HTML \rangle \\ H &\rightarrow HEAD BODY \mid HEAD \mid BODY \mid \lambda \\ HEAD &\rightarrow \langle HEAD \rangle HE \langle /HEAD \rangle \\ BODY &\rightarrow \langle BODY \rangle B \langle /BODY \rangle \\ HE &\rightarrow SCS \langle TITLE \rangle texto \langle /TITLE \rangle SCS \\ SCS &\rightarrow SC SCS \mid \lambda \\ SC &\rightarrow \langle SCRIPT \rangle TSC \langle /SCRIPT \rangle \\ TSC &\rightarrow \langle HTML \rangle TSC \mid \langle /HTML \rangle TSC \mid \\ &\quad \langle HEAD \rangle TSC \mid \langle /HEAD \rangle TSC \mid \\ &\quad \langle BODY \rangle TSC \mid \langle /BODY \rangle TSC \mid \\ &\quad \langle TITLE \rangle TSC \mid \langle /TITLE \rangle TSC \mid \\ &\quad \langle DIV \rangle TSC \mid \langle /DIV \rangle TSC \mid \\ &\quad \langle H1 \rangle TSC \mid \langle /H1 \rangle TSC \mid \\ &\quad \langle P \rangle TSC \mid \langle /P \rangle TSC \mid \\ &\quad \langle SCRIPT \rangle TSC \mid \langle BR \rangle TSC \mid texto TSC \mid \lambda \\ B &\rightarrow texto B \mid \\ &\quad \langle DIV \rangle B \langle /DIV \rangle B \mid \\ &\quad \langle H1 \rangle B \langle /H1 \rangle B \mid \\ &\quad \langle P \rangle B \langle /P \rangle B \mid \\ &\quad \langle BR \rangle B \mid \\ &\quad \lambda \end{aligned}$$

4 Arbol de derivación

A continuación presentamos el árbol de derivación para el siguiente ejemplo:

```
<HTML>
  <HEAD>
    <TITLE>Una pagina de ejemplo</TITLE>
    <SCRIPT>
      function unaFunc(){
        alert("esta funcion imprime un tag roto <TITLE>");
      }
    </SCRIPT>
    <SCRIPT></SCRIPT>
    <SCRIPT>alert("aca no aparece el cierre de SCRIPT")</SCRIPT>
  </HEAD>
  <BODY>
    <H1>Un heading</H1>
    <DIV>
      <P>Este texto es de prueba</P>
    </DIV>
    <BR>
    <P>Mas prueba</P>
  </BODY>
</HTML>
```

Debido al tamaño del ejemplo, el árbol de derivación tuvo que ser partido en dos partes. Nótese que del (único) nodo con el no terminal H cuelgan 2 hijos con el no-terminal $HEAD$, y el no terminal $BODY$. El primer subárbol muestra los descendientes de $HEAD$ y el segundo los de $BODY$. En cada árbol, el no terminal cuyo subárbol está suprimido se encuentra marcado en rojo, y con un subárbol genérico (triangulo) simbolizando que en realidad continúa. Por ejemplo:



En la figura *a.* se suprime el subárbol correspondiente al no terminal *BODY* y en el *b.* el correspondiente al no terminal *HEAD*.

