Parsování HTML představuje slangový výraz pro syntaktickou analýzu obsahu webové stránky. Lidově řečeno porcujeme zdrojový kód stránky a vyzobáváme potřebný obsah. Webovým vývojářům není neznámý pojem HTML DOM (Document Object Model). Ten je popsán konsorciem W3C (http://www.w3schools.com/jsref/dom_obj_document.asp) a umožňuje spravovat obsah stránky s pomocí javascriptu. Ačkoliv je možné brát inspiraci přímo z něj, v případě VBA se odvoláváme na starší knihovnu Microsoft HTML Object Library, v níž jsou některé vlastnosti definovány odlišně (outerHTML, innerText aj.). Každopádně výhodu mají ti, kteří se již potkali s vytvářením stránek v jazyce HTML a ovládají práci s tagy (elementy), jako je např.
body>, <div>, , , jejich atributy (vlastnostmi) a stylováním (CSS).

Pro účely testování jsem vytvořil stránku parsovani.html.

Zdrojový kód:

1	< !DOCTYPE html>
2	<html></html>
3	<head></head>
4	<title>Moje stránka</title>
5	<meta charset="utf-8"/>
6	<style></th></tr><tr><th>7</th><th>table {</th></tr><tr><th>8</th><th>border-collapse: collapse;</th></tr><tr><th>9</th><th>}</th></tr><tr><th>10</th><th>table, td, th {</th></tr><tr><th>11</th><th>border: 1px solid #CCCCCC;</th></tr><tr><th>12</th><th>}</th></tr><tr><th>13</th><th>.moje trida {</th></tr><tr><th>14</th><th>color: #0066CC;</th></tr><tr><th>15</th><th>}</th></tr><tr><th>16</th><th></style>
17	
18	<body></body>
19	<input name="muj_nazev" type="text" value="Excel 2010"/>
20	Element P s atributem
21	id="muj_identifikator". <b p>
22	< div class="moje_trida">
23	Element P v prvním elementu DIV s atributem class="moje_trida"
24	(index 0).
25	
26	< div class="moje_trida">
27	Element P ve druhém elementu DIV s atributem class="moje_trida"
28	(index 1).
29	
30	
31	
32	Křížek
33	123,45
34	20.11.2015
35	
36	
3/	Bydzovsky
38	6/8,90
39	<ta>1.6.2016</ta>
40	
41	
42	<img alt="<br"/> are ČT1 beight 02 are at1 is glussidate 155 /beight 1
43	"Logo CIII" neight="83" SrC="Ct1.jpg" Width="155"/> bref. "bttp://www.gooleatelouing.go/et2/" toward. "blank", since alt
44	<img alt="</th"/>
40 76	Logo CTZ meignt= δS SrC= CtZ.jpg width="100"/>
40	
4/	

A nyní už si pojďme obsah stránky rozebrat programově. Kód je taktéž uveden v příloze a doporučuji jej krokovat a studovat v oknech Immediate a Locals.

'Tools / References / Microsoft HTML Object Library

Dim objMSHTML As New HTMLDocument Dim objDocument As HTMLDocument

Sub ParsovaniHTML()

Dim objimages As IHTMLElementCollection Dim objElements As IHTMLElementCollection Dim objElements As IHTMLElementCollection Dim objTags Sa IHTMLElementCollection Dim objTagsTagName As IHTMLElementCollection Dim objTagsSas As IHTMLElementCollection Dim objTagsName As IHTMLElementCollection

Dim objImage As IHTMLImgElement Dim objLink As IHTMLAnchorElement

Dim objTagClass As IHTMLElement Dim objTagName As IHTMLElement Dim objTagId As IHTMLElement

'přesnější typy vyplývající z testování 'Dim objTagClass As IHTMLDivElement 'Dim objTagName As IHTMLInputElement 'Dim objTagId As IHTMLParaElement

Dim objHTML As IHTMLHtmlElement Dim objHead As IHTMLHeadElement Dim objBody As IHTMLBodyElement

Dim strURL As String Dim strTitulek As String Dim strHTML As String Dim strHead As String Dim strBedy As String

'adresa stránky strURL = "http://proexcel.cz/test/parsovani.html"

'element ... tag

'dokument Set objDocument = objMSHTML.createDocumentFromUrl(strURL, vbNullString) 'čekání na stažení

While objDocument.readyState <> "complete" DoEvents Wend

'<mark>titulek stránky</mark> strTitulek = objDocument.title 'objekt HTMI (element html)

'objekt HTML (element html) Set objHTML = objDocument.documentElement strHTML = objHTML.outerHTML 'blavička (element head)

'hlavička (element head) Set objHead = objDocument.head strHead = objHead.outerHTML

'obsah stránky (element body) Set objBody = objDocument.body strBody = objBody.outerHTML

'kolekce obrázků (elementy img) Set objimages = objDocument.images

For Each objimage In objimages Debug Print objimage.outerHTML Debug Print objimage.getAttribute("href") Next

'kolekce hypertextových odkazů (elementy a) Set objLinks = objDocument.links

For Each objLink In objLinks Debug.Print objLink.outerHTML Debug.Print objLink.innerHTML Debug.Print objLink.getAttribute("href") Next

'varianta 1 pro tagy

'kolekce elementů Set objElements = objDocument.all

'kolekce elementů s požadovaným názvem (p) Set objTags = objElements.tags("p") 'varianta 2

'kolekce elementů s požadovaným názvem (p) Set objTagsTagName = objDocument.getElementsByTagName("p")

'element s atributem id (id="muj_identifikator") "ID by mělo být v dokumentu jedinečné Set objTagId = objDocument.getElementById("muj_identifikator")

'typ nalezeného elementu

strElement = objTagId.tagName

'získání barvy atributu style nalezeného elementu (style="color: ...") '##0066 strColor = objTagld.style.Color

'kolekce elementů s požadovaným atributem class (class="moje_trida") Set objTagsClass = objDocument.getElementsByClassName("moje_trida")

For Each objTagClass In objTagSClass Debug.Print objTagClass.tagName Debug.Print objTagClass.outerHTML Debug.Print objTagClass.innerHTML Next

'kolekce elementů (zpravidla elementy input) 's požadovaným atributem name (name="hledaný řetězec") Set objTagsName = objDocument.getElementsByName("muj_nazev")

For Each objTagName In objTagsName Debug.Print objTagName.tagName Debug.Print objTagName.outerHTML Debug.Print objTagName.getAttribute("value") Next

'odstranění z paměti Set objDocument = Nothing Set objMSHTML = Nothing

End Sub

Řádky VBA jsem se snažil komentovat a na tomto místě jen upřesním pojmy innerHTML, innerText a outerHTML.

Příklad <div>nějaký text</div>

```
<div>nějaký text</div> ... vlastnost outerHTML pro element <div>nějaký text ... vlastnost innerHTML pro element <div>
nějaký text ... vlastnost innerTEXT pro element
```

Pozn. Pokud se chcete odkazovat na členy kolekcí indexem, pak vězte, že číslování začíná nulou.

Přirozeně se sluší na tomto místě ukázat způsob, jak z dané stránky převzít tabulku do listu Excelu (ačkoliv prosté HTML tabulky je lepší načítat prostřednictvím karty Data / Z webu).

1	Sub ParsovaniTabulkyHTML()
2	Dim obiMSHTML As New HTML Decument
4	Dim objocument As HTMI Document
5	Din objectment Re InfileDocument
6	Dim objTagsRow As IHTMLElementCollection
7	Dim objTagsCell As IHTMLElementCollection
8	
9	Dim objTagRow As IHTMLTableRow
10	Dim objTagCell As IHTMLTableCell
11	
12	Dim strukl As String
14	'adresa stránky
15	strURL = "http://proexcel.cz/test/parsovani.html"
16	
17	'dokument
18	Set objDocument = objMSHTML.createDocumentFromUrl(strURL, vbNullString)
19	
20	'čekání na stažení
21	while objDocument.readyState <> "complete"
22	Vend
23	Welld
25	'všechny řádky tabulky
26	Set objTagsRow = objDocument.getElementsByTagName("tr")
27	
28	'pro každý řádek
29	For Each objTagRow In objTagsRow
30	
31	vsecny bunky radku
32	Set $objragscen = objragkow.getelementsbyragName("to")$
34	'čítač pro řádky
35	i = i + 1
36	
37	For Each objTagCell In objTagsCell
38	
39	'čítač pro sloupce
40	J = J + I
41	zánis do huněk listu
43	
44	Case 1
45	'text
46	Cells(i, j).Value = objTagCell.innerText
47	Case 2
48	'desetinné číslo
49	Cells(I, J).Value = CDbl(objfagCell.innerText)
50	Case 3
52	Cells(i, i) Value = CDate(obiTagCell innerText)
53	End Select
54	
55	Next objTagCell
56	
57	'reset čítače pro sloupce
58	$\mathbf{J} = 0$
59	Next obiTagBow
61	NEAL OUT AGROW
62	'odstranění z paměti
63	Set objDocument = Nothing
64	Set objMSHTML = Nothing
65	
66	End Sub

×

Tabulka převedená z HTML stránky

HTML stránky by do jisté míry měly dodržovat hierarchii objektů a jejich vnořování do sebe. V praxi tomu tak často není a jejich obsah bývá uspořádán laxně, na rozdíl třeba od XML. Je to jeden z důvodů, proč i já jsem v daném tématu nevyužil skutečnosti, že tagy (elementy) představují jakési "nody" ve stromové struktuře, kdy uvažujeme vazby rodič (parent) – dítě (child), případně děti (children).

Parsování HTML stránek nepatří k technikám, za které bychom se mohli plácat po ramenou. Pokud máme možnost, vždy sáhneme po přímém přístupu k datům do databáze. Klíčové je slovíčko "pokud". Až příliš dobře se pamatuji na nutnost zpracovat 60 000 webových stránek z nejmenovaného webu státní správy jen proto, že webová aplikace padala pod deseti minutách nastavování parametrů (bez možnosti uložení). Poměrně solidně se s HTML kódem umí vypořádat i regulární výrazy.

Ke zpracování webových stránek a jejich obsahu se opět někdy vrátíme.

Příloha

excel_parsovani_html.zip