

Pattern for Data Cleanup

Copyright © 2017 Chandra Lingam

Write a pattern to cleanup and prepare the HTML documents for CSV Conversion. Here are two sample texts provided.

If you notice the below text, we need to extract values inside "tr", "td", "th" elements. However, HTML is very forgiving and flexible standard and allows liberal use of white spaces and newlines.

To simplify parsing work, we want to first convert the incoming HTML document and replace one or more lines, tabs, and spaces with a single space.

The pattern you write should match: spaces, new line character tabs, and HTML code for space (view labresults.html in a text editor to see the use of html code).

When you replace the matching text with a single space, the net effect is HTML content is converted to a single line text with single space replacing all newlines, tabs, multiple spaces and html code for space

```
<tr ID="problem59" styleCode="normRow">
    <td ID="problem59name">At high risk for heart block</td>
    <td>02/10/2015</td>
    <td>06/15/2015</td>
</tr>
```

```
<tr><th align="left">Hemoccult
```

```
</th><td></td><td> </td><td>negative</td><td> </td><td> </td></tr>
```