



An example webpage

- For instance, here's the source code of a very basic webpage.

```
In [14]: 1 html_string = '''
2 <html>
3     <body>
4         <div id="content">
5             <h1>Heading here</h1>
6             <p>My First paragraph</p>
7             <p>My <em>second</em> paragraph</p>
8             <hr>
9         </div>
10        <div id="nav">
11            <ul>
12                <li>item 1</li>
13                <li>item 2</li>
14                <li>item 3</li>
15            </ul>
16        </div>
17    </body>
18 </html>'''
```

- Here's what that webpage actually looks like:

```
In [15]: 1 HTML(html_string)
```

Out[15]:

Heading here

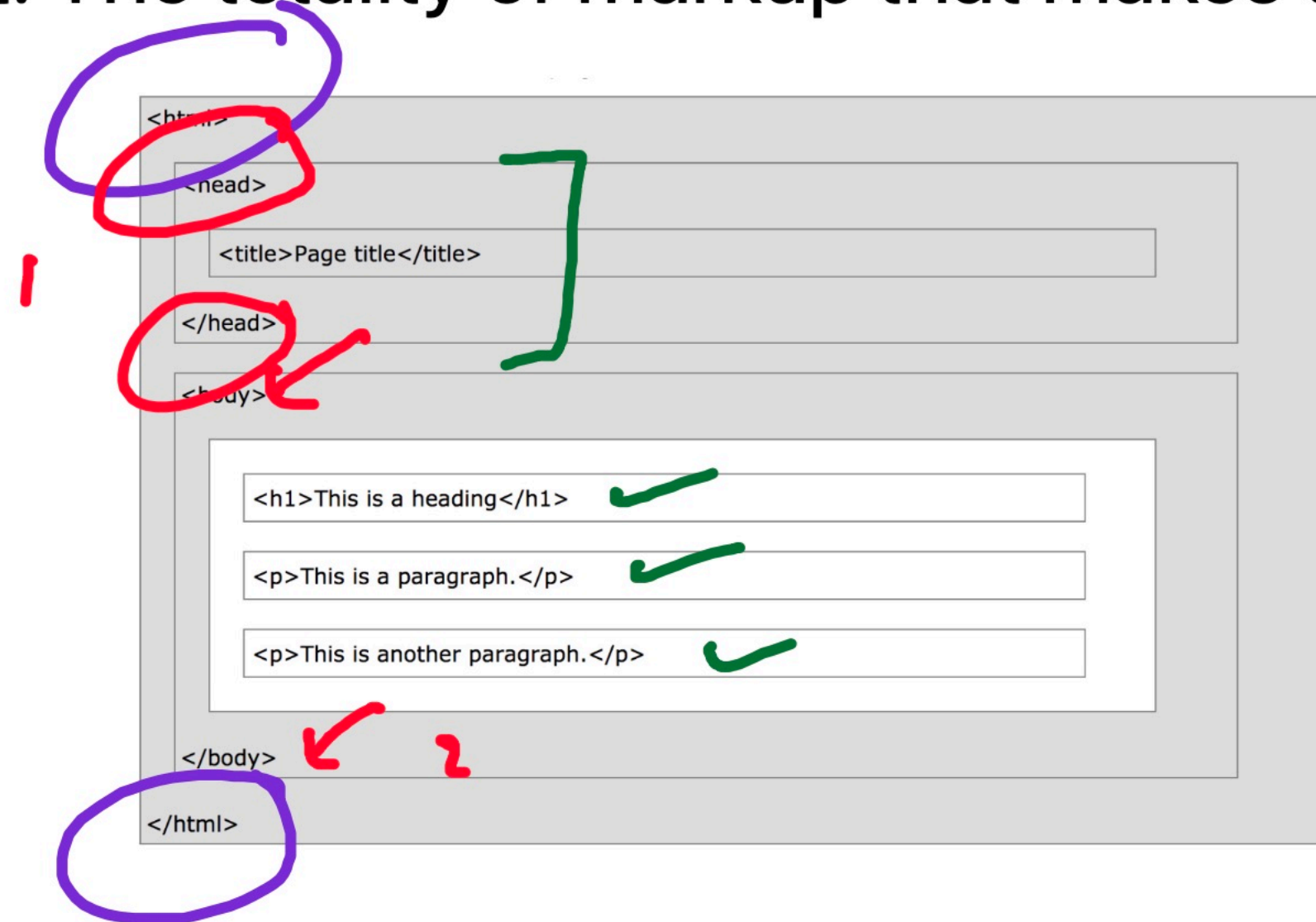
My First paragraph

My *second* paragraph

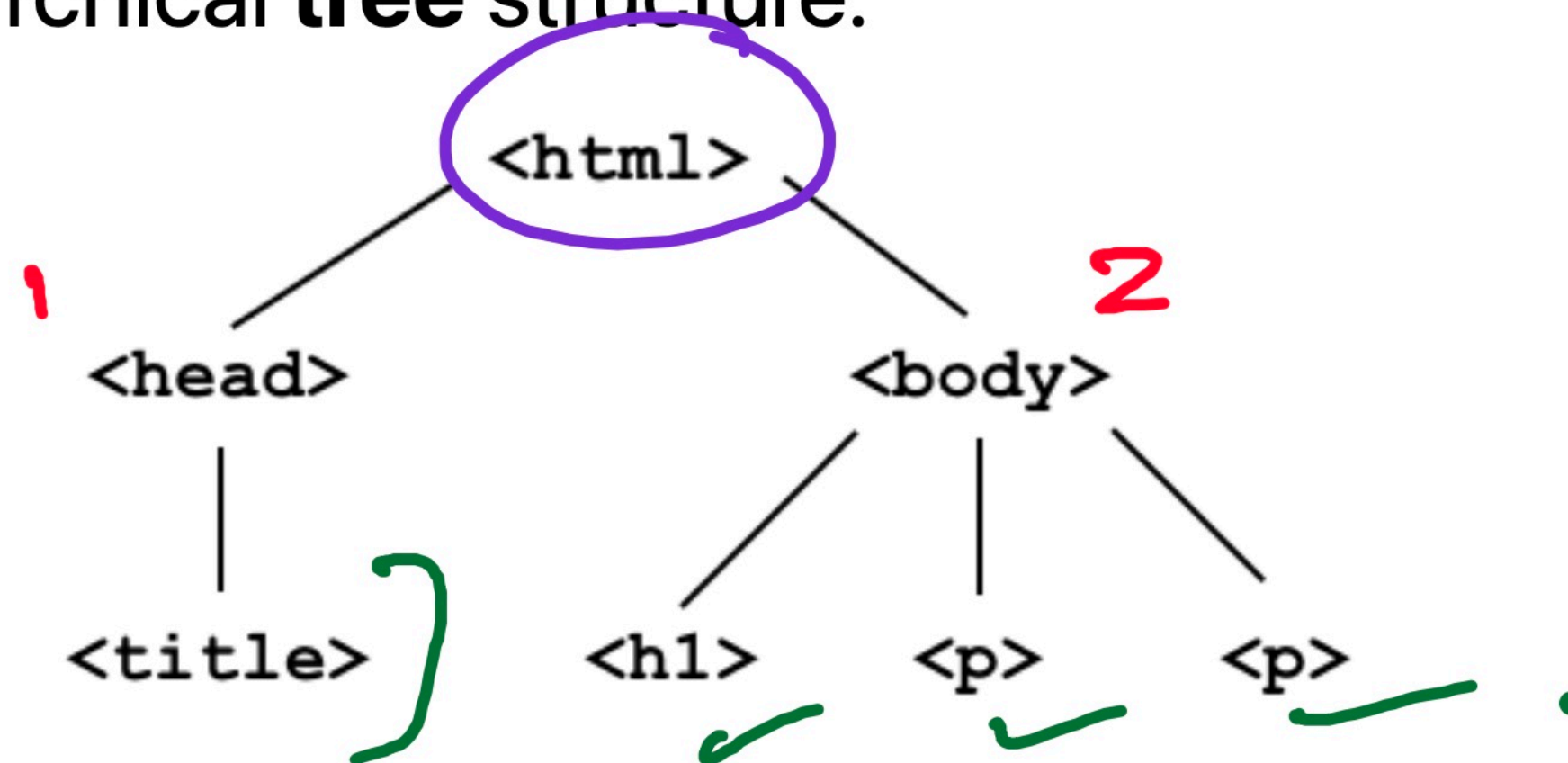
- item 1
- item 2
- item 3



- **HTML document:** The totality of markup that makes up a webpage.



- **Document Object Model (DOM):** The internal representation of an HTML document as a hierarchical **tree** structure.



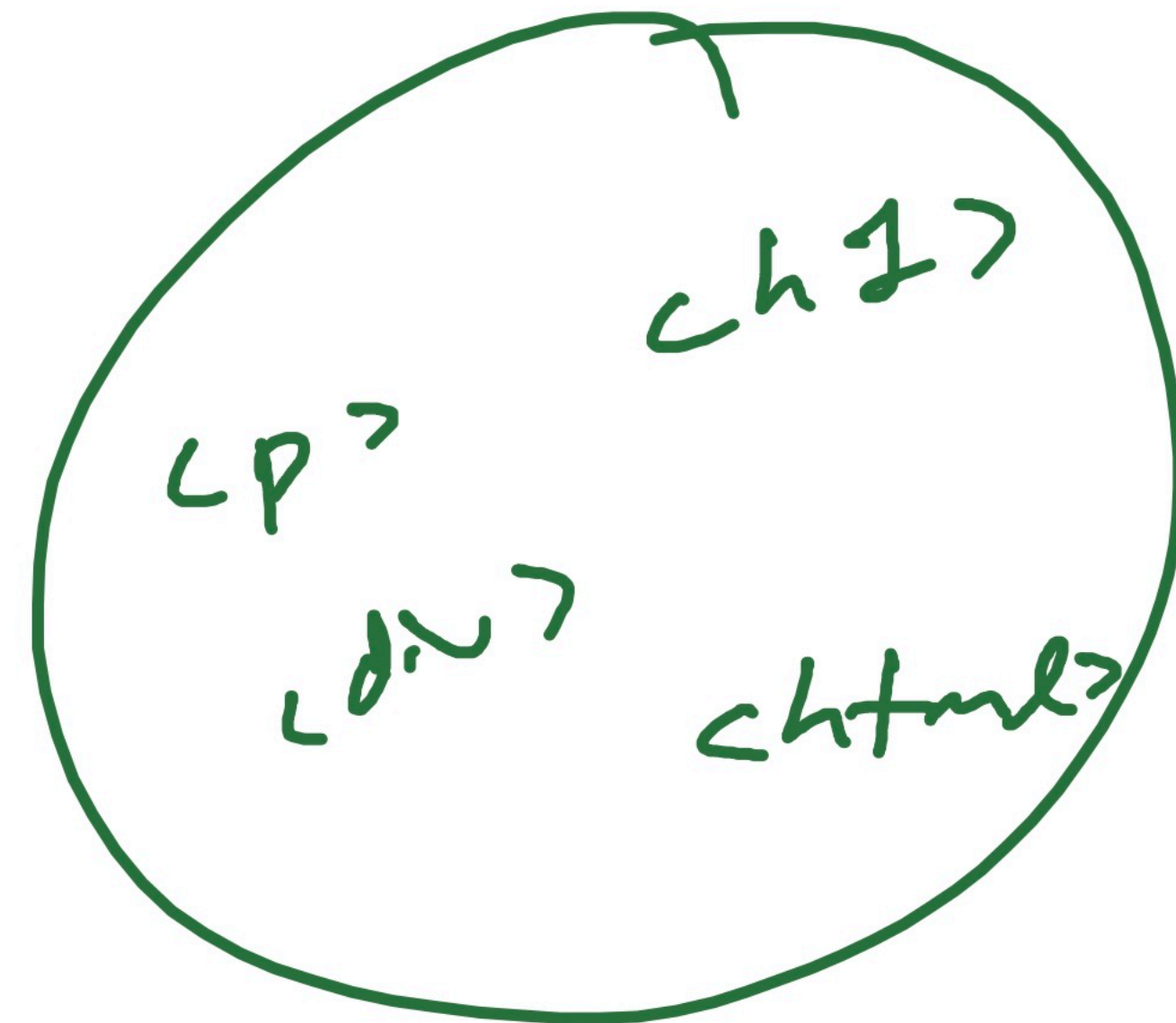
Beautiful Soup

- Beautiful Soup 4 is a Python HTML parser.

Remember, **parse** means to "extract meaning from a sequence of symbols".

- **Warning:** Beautiful Soup 4 and Beautiful Soup 3 work differently, so make sure you are using and looking at documentation for Beautiful Soup 4.

Rest assured, the `pds` conda environment already has Beautiful Soup 4 installed.




```

6  ``html
7  <head>
8    <title>3*Canada-2022-06-04</title>
9  </head>
10 <body>
11   <h1>Spotify Top 3 - Canada</h1>
12   <table>
13     <tr class='heading'>
14       <th>Rank</th>
15       <th>Artist(s)</th>
16       <th>Song</th>
17     </tr>
18     <tr class=1>
19       <td>1</td>
20       <td>Harry Styles</td>
21       <td>As It Was</td>
22     </tr>
23     <tr class=2>
24       <td>2</td>
25       <td>Jack Harlow</td>
26       <td>First Class</td>
27     </tr>
28     <tr class=3>
29       <td>3</td>
30       <td>Kendrick Lamar</td>
31       <td>N95</td>
32     </tr>
33   </table>
34 </body>
35 ``

```

14 leafs

- **Part 1:** How many leaf nodes are there in the DOM tree of the previous document — that is, how many nodes have no children?

```

8     <title>3*Canada-2022-06-04</title>
9 </head>
10 <body>
11     <h1>Spotify Top 3 - Canada</h1>
12     <table>
13         <tr class='heading'>
14             <th>Rank</th>
15             <th>Artist(s)</th>
16             <th>Song</th>
17         </tr>
18         <tr class=1>
19             <td>1</td> ✓
20             <td>Harry Styles</td> ✓
21             <td>As It Was</td> ✓
22         </tr>
23         <tr class=2>
24             <td>2</td> ✓
25             <td>Jack Harlow</td> ✓
26             <td>First Class</td> ✓
27         </tr>
28         <tr class=3>
29             <td>3</td> ✓
30             <td>Kendrick Lamar</td> ✓
31             <td>N95</td> ✓
32         </tr>
33     </table>
34 </body>
35

```

- **Part 1:** How many leaf nodes are there in the DOM tree of the previous document — that is, how many nodes have no children?

- **Part 2:** What does the following line of code evaluate to?

```
len(soup.find_all("td"))
```



```

1 <div class="alert alert-success" markdown="1">
2   <h3>Activity</h3>
3
4 Consider the following HTML document, which represents a webpage containing the top few songs with the most streams on Spotify today in Canada.
5
6 ```html
7 <head>
8   <title>3*Canada-2022-06-04</title>
9 </head>
10 <body>
11   <h1>Spotify Top 3 - Canada</h1>
12   <table>
13     <tr class='heading'>
14       <th>Rank</th>
15       <th>Artist(s)</th>
16       <th>Song</th>
17     </tr>
18     <tr class=1>
19       <td>1</td>
20       <td>Harry Styles</td>
21       <td>As It Was</td>
22     </tr>
23     <tr class=2>
24       <td>2</td>
25       <td>Jack Harlow</td>
26       <td>First Class</td>
27     </tr>
28     <tr class=3>
29       <td>3</td>
30       <td>Kendrick Lamar</td>
31       <td>N95</td>
32     </tr>
33   </table>
34 </body>
35 ```

```

- **Part 1:** How many leaf nodes are there in the DOM tree of the previous document — that is, how many nodes have no children?
- **Part 2:** What does the following line of code evaluate to?

```
len(soup.find_all("td"))
```
- **Part 3:** What does the following line of code evaluate to?

```
soup.find("tr").get("class")
```