**11.4.1**

## Gaining a Deeper Understanding of a Website Structure

**Now** that she knows basic HTML and CSS, Robin wants to find specific data. But on a webpage with lots of content, finding the HTML element that contains the data she wants can be overwhelming. So, she'll use Chrome Developer Tools (DevTools). With DevTools, developers can review the structure of any webpage. And not only that, but it also has a search function. This will help make sense of the tags and elements holding the data that Robin is seeking.

Specifically, Robin needs to identify elements of a webpage so that she can extract data from the News – NASA Mars Exploration ⬈ (https://mars.nasa.gov/news/? page=0&per_page=40&order=publish_date+desc%2Ccreated_at+desc&search=&category=19%2C165%2C18 4%2C204&blank_scope=Latest) website. Let's follow along with her.

You've already learned how to identify and create basic HTML code in preparation for web scraping. You're now ready to learn more advanced skills to further your knowledge of how a website is structured. You'll follow the same process for each webpage that you want to scrape: visit the page, identify the data, then go through the HTML code to pinpoint its location on the webpage.

In this section, you'll learn how to use CSS selectors to efficiently find elements in Chrome DevTools.
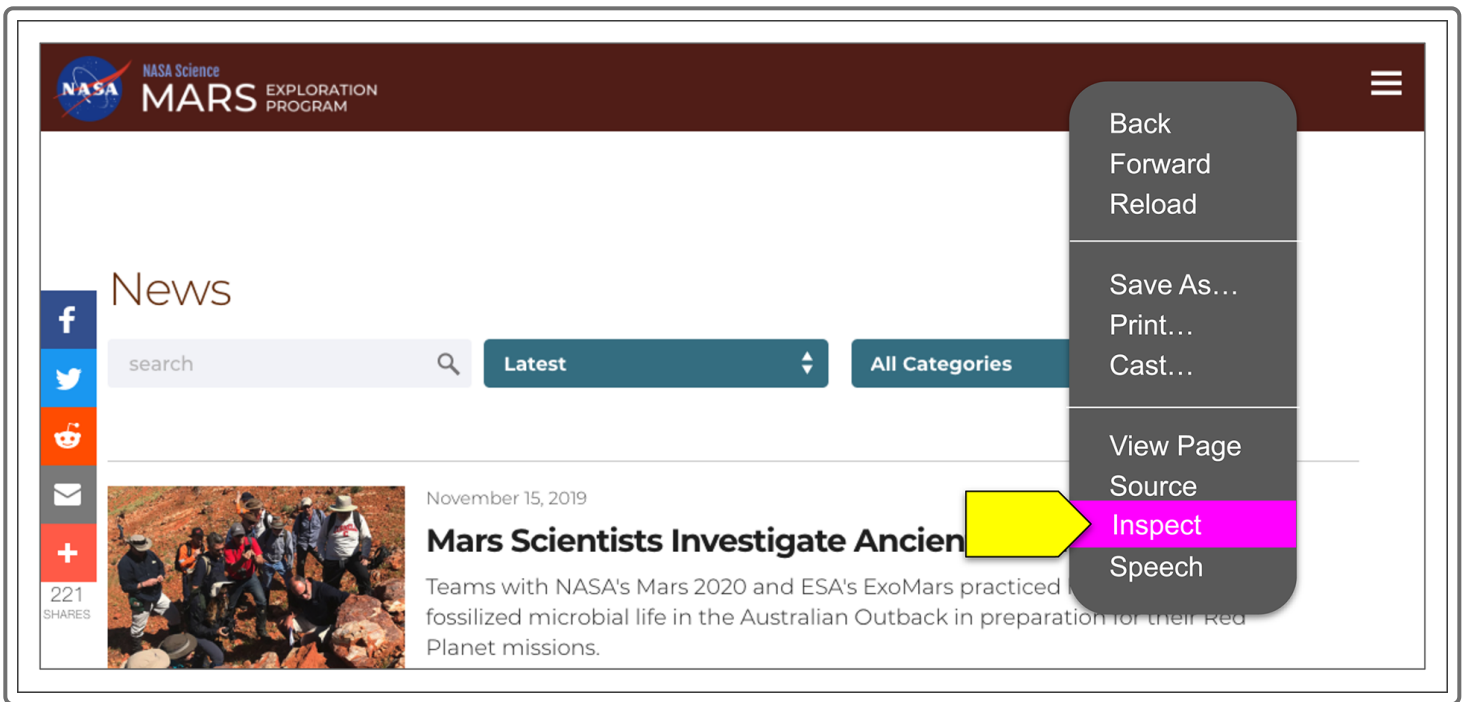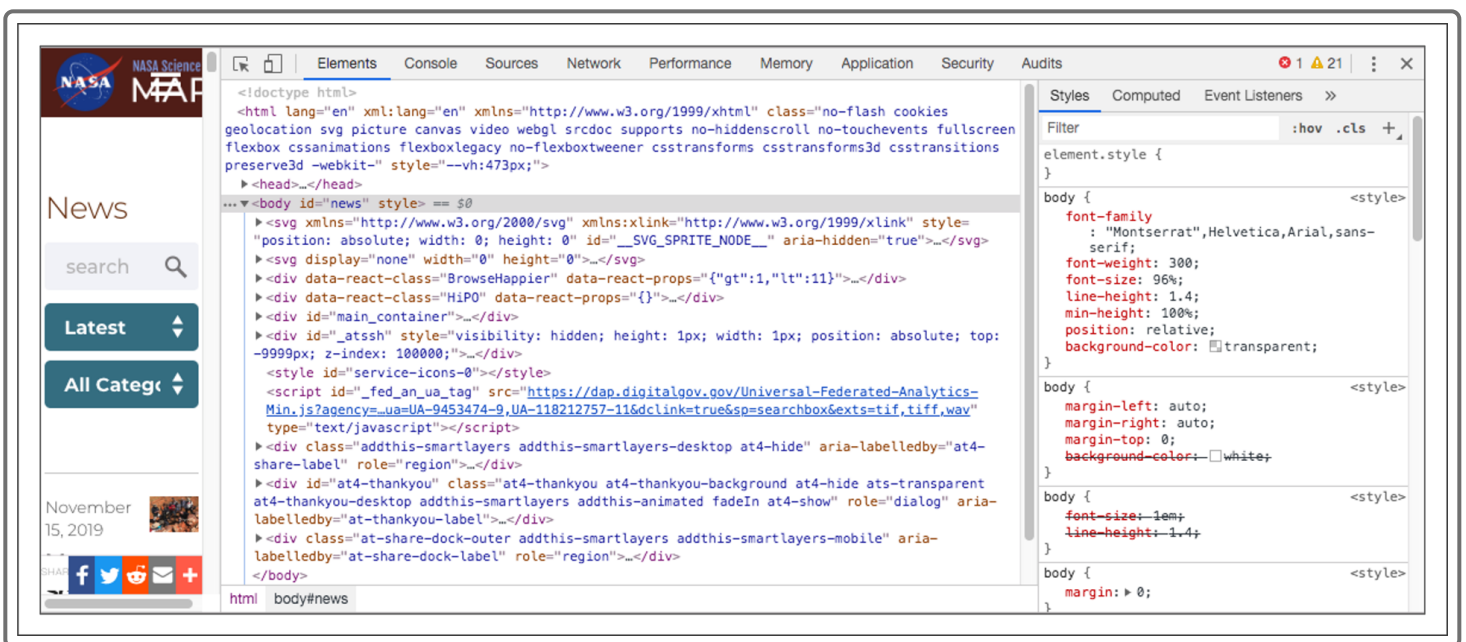
## Using Chrome DevTools

As you scrape websites, you'll often find that the HTML element you want is contained within another element—or even buried inside layers of elements. DevTools makes it easier to identify such an element and its CSS selectors. We can then use those selectors to scrape the data from the element.

Start by opening the News – NASA Mars Exploration ⬈ (https://mars.nasa.gov/news/? page=0&per_page=40&order=publish_date+desc%2Ccreated_at+desc&search=&category=19%2C165%2C184%2C 204&blank_scope=Latest) website in a new browser window. We can initially observe article titles and a sentence that describes each article.

Next, open DevTools by right-clicking anywhere on the page and then clicking Inspect.

A new panel opens. This panel is docked to the webpage but has a different job than the webpage.



As we can observe in the panel, this site has a lot going on. What's the overall structure of the site? Well, the line that begins with `<html lang="en">` should seem familiar. And, so should the `<head />` and `<body />` tags. But, what's all the other stuff? And, what's the stuff that appears inside the familiar tags?

Drag the HTML tags to the correct positions in the code.

```
<[          ] id="news" style="">
  <[          ] display="none" width="0" height="0">
    <[          ] id="circle_plus" height="30" viewBox="0 0 30 30" width="30">
      <[          ] fill-rule="evenodd" transform="translate(1 1)">
        <[          ] cx="14" cy="14" ...><[          ]>
        <[          ] class="the_plus" d="m18.856 ...><[          ]>
      <[          ]>
    <[          ]>
  <[          ]>
```
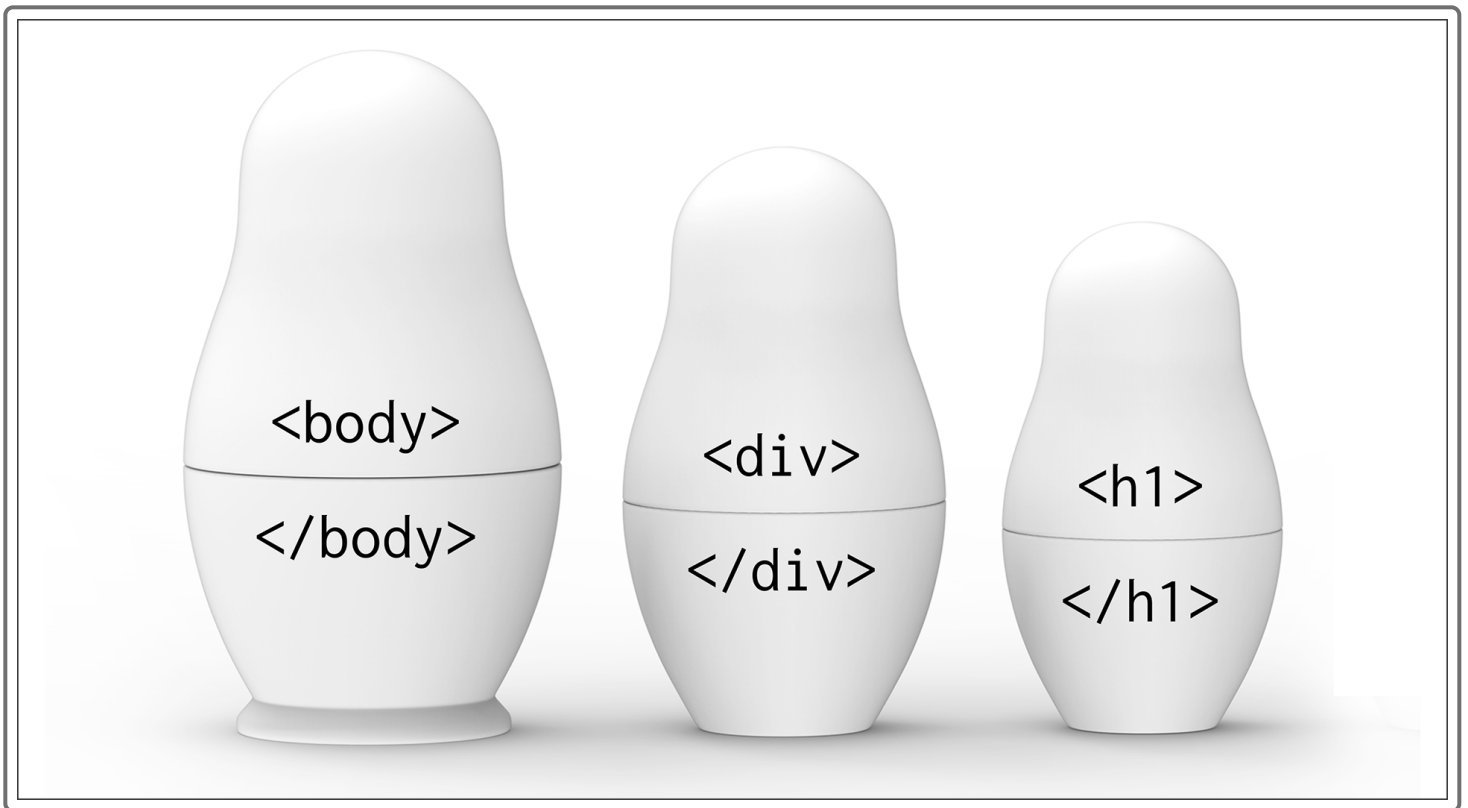
| :: g | :: circle | :: /circle | :: /symbol | :: path | :: svg | :: body | :: /body |
|------|-----------|-----------|-----------|---------|--------|---------|----------|

| :: /svg | :: /path | :: /g | :: symbol |
|---------|----------|-------|-----------|

Check Answer

Finish ▶

Let's break down the structure of this site a bit. Remember that an HTML element can contain other elements. For example, the `<body />` tag is a container of every element that appears on the webpage, such as the header and paragraph elements. And, other containers can also exist inside that `<body />` tag. These containers are **nested**, much like in a nesting doll. That is, one element can contain another element which, in turn, can contain another
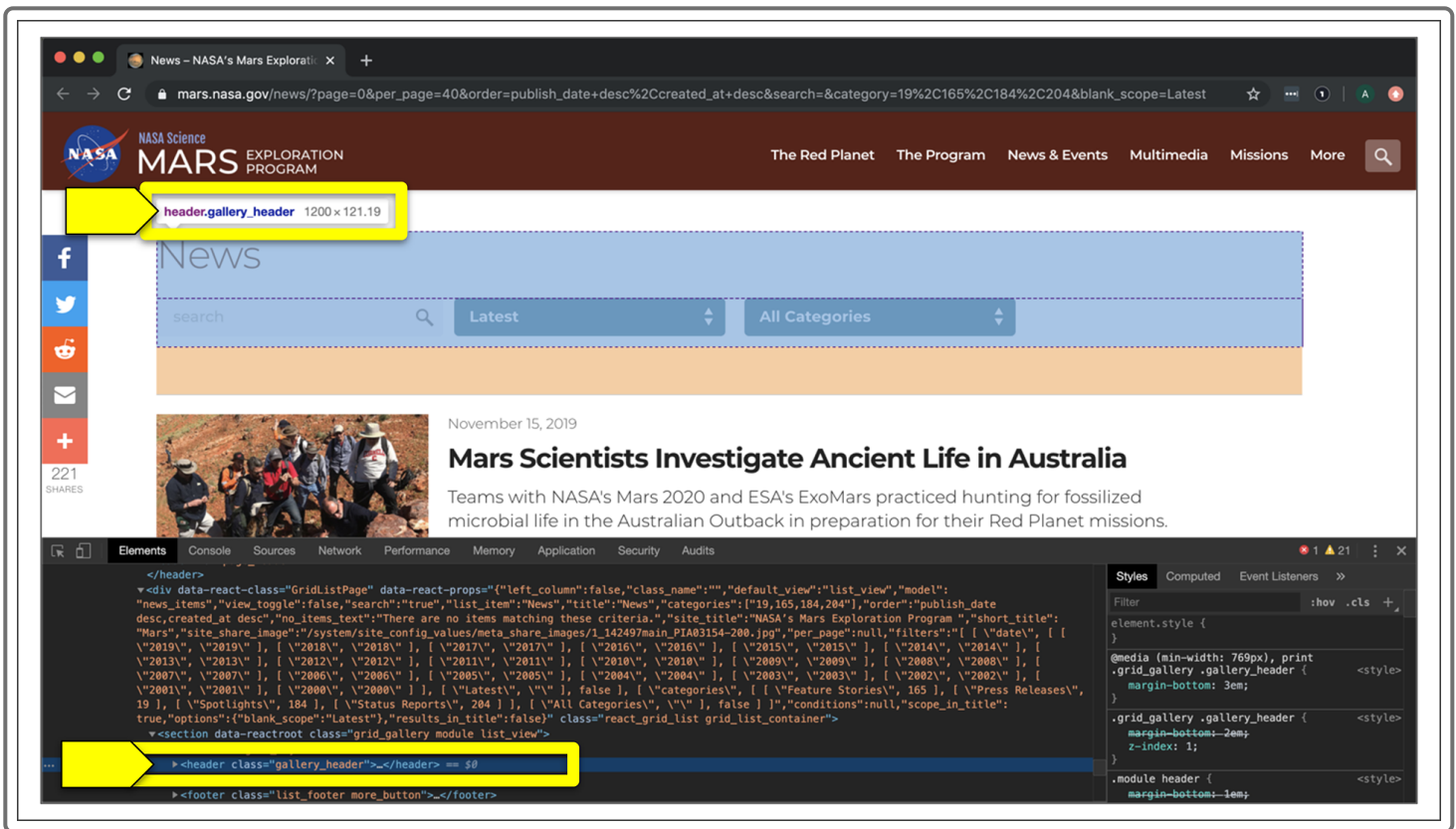
element, and so on. Furthermore, multiple levels of nesting can exist, depending on how elaborate the website is. In the case of this website (and most websites), the other containers inside the body consist of `<div />` tags.

The following image shows that each container is nested within another. In this case, an `<h1 />` element is nested within a `<div >` container which, in turn, is nested within the `<body />` container:



DevTools has another feature. Specifically, when we hover over any part of the code, the connected visual gets highlighted on the webpage. This helps by showing us which piece of code is tied to which feature of the webpage.

This webpage contains lots of custom code. So instead of scrolling through all of it to find a certain element, we'll just search for it in DevTools. To do so, press Ctrl+F (on Windows) or Command+F (on macOS) to display the Find box. In the box, enter "gallery_header" (without the quotation marks), and then press Enter. It may return more than one result, so you may have to press Enter to advance until the "header class="gallery_header" line is selected, and then hover over it. Make sure that the "header class="gallery_header" line is selected, and then hover over it. The header section of the page, which consists of the title and its container element, gets highlighted.

Next, hover over the following line of code, which is `<h2 class="module_title">News</h2>`. (If the header doesn't display the nested contents, click its arrow to expand it.)

The highlighted portion of the webpage becomes smaller. That's because we hovered over an element that's nested inside a container instead of hovering over the full container.

This is a splendid way to pinpoint where on the website we want our web-scraping code to extract data from. We can't tell the code to just extract a division or a header, though. That's because many of these might exist on the webpage, when we want only one. And, that's where the `class` and `id` attributes come into play.

## © 2022 edX Boot Camps LLC