# Chapter 9

# Using Development Tools to Examine Webpages

**Skills you will learn:** For this tutorial, we will use the developer tools in Firefox. However, these are quite similar to the developer tools found in just about every modern web browser, including Google's Chrome and Microsoft's Edge.

We aren't covering everything Firefox's developer tools can do, but instead are focusing on aspects that are of value for web scraping and simple development.

Note that the prior version of this tutorial focused on the Firebug plugin, which has since been integrated into Firefox's toolset.

**Summary:** Hidden behind the public face of modern web browsers is a suite of tools that gives you extraordinary power to inspect the source code of web pages, and observe network requests and responses. These tools are useful both in web scraping, discussed in chapter 9, and web development, discussed in chapter 10.

**Introduction:**

All web browsers allow the user to look at the source HTML for the page, and this is perfectly useful for looking at the overall HTML structure of a page, for seeing tags in their full context, for searching for elements, and so on. But sometimes you want to drill down and look more closely at specific elements on the page, or you want to be able to see the network traffic being passed to and from the remote web server. For these tasks, development tools can be the perfect solution.

**Using the tools**

The developer tools are an integrated part of Firefox.

To start using them, got to Tools>Web Developer>Toggle Tools. We'll use the Mac version in this tutorial, but the functionality is essentially identical on a Windows PC.

**Examining Elements in a webpage**

We're going to have a look at the source code for the Public Works and Government Services Canada (now renamed Public Services and Procurement Canada) proactive disclosure page that we first saw in chapter 9 of *The Data Journalist*. While the government is moving the proactive disclosure information to its open data portal, the page you saw in Chapter 9 is archived at this address: https://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=L;Sort=0;PF=CL201516Q1.txt

If you simply open the page source (Tools>Web Developer>Page Source), you will see all of the HTML for the page.

```html
<h1 id="wb-cont">2015-2016 - 1<sup>st</sup> Quarter (April to June)</h1><p class="note">N.B.: The contract date represents the date that the contrac
<caption class="font-small align-left">
<em>List of contracts for the trimester, which includes the date, vendor, description and value</em>
</caption>
    <tr >
        <th class="secondary width20" scope="col"><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=L;Sort=0;PF=CL201516Q1.tx
        <th class="secondary" scope="col"><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=L;Sort=1;PF=CL201516Q1.txt" title
        <th class="secondary" scope="col"><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=L;Sort=2;PF=CL201516Q1.txt" title
        <th class="secondary" scope="col"><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=L;Sort=3;PF=CL201516Q1.txt" title
    </tr>
    <tr>
        <td>2008-04-22</td>
        <td><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=228" title="IGF VIGILANCE INC.;20
        <td>420 - Engineering Services not Elsewhere Specified</td>
        <td class="alignRight">$500,695.10</td>
    </tr>
    <tr>
        <td>2008-04-22</td>
        <td><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=432" title="THE ARCOP GROUP &nbsp
        <td>423 - Engineering Consultants - Other</td>
        <td class="alignRight">$67,145,777.15</td>
    </tr>
    <tr>
        <td>2008-04-22</td>
        <td><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=1698" title="CAPITAL ELEVATOR LTI
        <td>859 - Other Business Services not Elsewhere Specified</td>
        <td class="alignRight">$73,651.02</td>
    </tr>
    <tr>
        <td>2008-09-04</td>
        <td><a href="http://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=266" title="NORR LIMITED;2008-09-
        <td>421 - Architectural Services</td>
        <td class="alignRight">$40,100,544.85</td>
    </tr>
    <tr>
        <td>2009-02-28</td>
```

Looking at the HTML code this way is useful way to put it all in context, especially when you already know what elements you want to look at or when you are trying to figure out how you will instruct a Python library such as Beautiful Soup to pinpoint the elements you want (Chapter 9 introduces Beautiful Soup).

But when you are first exploring a page to see how it works, the element inspector can make the job a lot easier.

If you open Developer Tools while the page is open, you see a panel like the one below appears at the bottom of the browser screen.

**2015-2016 - 1st Quarter (April to June)**

N.B.: The contract date represents the date that the contract is recorded in the departmental financial system.

*List of contracts for the trimester, which includes the date, vendor, description and value*

| Contract Date | Vendor Name | Description of Work | Contract Value |
|---|---|---|---|
| 2008-04-22 | IGF VIGILANCE INC. | 420 - Engineering Services not Elsewhere Specified | $500,695.10 |
| 2008-04-22 | THE ARCOP GROUP / GERSOVITZ MOSS | 423 - Engineering Consultants - Other | $67,145,777.15 |
| 2008-04-22 | CAPITAL ELEVATOR LTD. | 859 - Other Business Services not Elsewhere Specified | $73,651.02 |
| 2008-09-04 | NORR LIMITED | 421 - Architectural Services | $40,100,544.85 |
| 2009-02-28 | R&R AUTOMATION INC. | 859 - Other Business Services not Elsewhere Specified | $1,077,864.27 |
| 2009-03-25 | CAPITAL ELEVATOR LTD. | 859 - Other Business Services not Elsewhere | $108,703.30 |

If you look at the top left corner of the panel, you can see the element inspector icon.



Clicking on the icon makes the inspector active. With it we can drill down into the HTML for any part of the page that we want to examine. Simply click on a part of the webpage in the upper window containing the web page as rendered by Firefox, and the inspector will show you the underlying HTML and CSS code. We'll use the inspector to click on the HTML table value IGF Vigilance Inc.

We then see the HTML for that part of the page, in the lower panel.



The small expansion arrows to the left of individual HTML elements indicate that more detail is available to be viewed. If we click on that, we see more code, and yet another expansion arrow indicating we can drill down even further.

```
        <td>2008-04-22</td>
      ▼ <td>
        ▶ <a href="https://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=228"
          title="IGF VIGILANCE INC.;2008-04-22">⊡</a> ev
        </td>
```

So we'll click on that to reveal the last layer of detail.

```
      ▼ <td>
        ▼ <a href="https://www.tpsgc-pwgsc.gc.ca/cgi-bin/proactive/cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q1.txt;LN=228"
          title="IGF VIGILANCE INC.;2008-04-22"> ev
            <span class="wb-hide">228</span>
            IGF VIGILANCE INC.
          </a>
        </td>
```

We can now see all of the HTML code associated with the IGF Vigilance entry. First off, it is contained in a <td> or table cell tag, which itself is enclosed in a <tr> or table row tag, which in turn is enclosed in a <tbody> or table body tag, and so in, in the hierarchy of tags that makes up the webpage.

Within the <td> tag is a <a> or hyperlink tag. The title attribute that follows causes its contents to appear as a tooltip when a user hovers the mouse over the link. This is then followed by the HREF attribute, which indicates the destination of the link, in this case another Public Words web page that provides detail on this particular contract. This is followed by a <span> element that has a CSS class that hides it on the page, and finally the actual text seen by the browser user. This is followed by the closing </a> tag for the hyperlink.

If we were scraping data from the page, we would probably be most interested in the text for the name of the contract vendor, IGF Vigilance Inc.  By using the element inspector, we have identified that we would be looking for the visible display text within the hyperlink within the <td> tag. We can use a parser such as Beautiful Soup, or a regular expression, to isolate the desired text element, and grab it for further processing (Chapter 9 walks through how to scrape this page using a copy of the data saved so it will remain available).

**Using the Network monitor tab**

As discussed in chapter 9, when you enter a web address in a browser, click on a link, or when a script running in the page requests an external resource, a request is sent to a remote server, somewhere else on the Internet, and a response comes back. You can use the Network panel in the developer tools to examine these requests and responses, so as to better understand what the browser is asking for, and what is being sent back in response.

We'll begin with a simple example, what happens when you click on one of the details links in the Proactive Disclosure page we've been looking at.

To start, open developer tools and click on the network tab.

Here's a look at the network panel:

**2015-2016 - 1ˢᵗ Quarter (April to June)**

Disclosure of Contracts

Disclosure of Travel and Hospitality Expenses

Disclosure of Annual Expenditures for Travel, Hospitality and Conferences

Disclosure of Position Reclassifications

Disclosure of Grants and Contributions

Disclosure of wrongdoings in the workplace

N.B.: The contract date represents the date that the contract is recorded in the departmental financial system.

*List of contracts for the trimester, which includes the date, vendor, description and value*

| Contract Date | Vendor Name | Description of Work | Contract Value |
|---|---|---|---|
| 2008-04-22 | IGF VIGILANCE INC. | 420 - Engineering Services not Elsewhere Specified | $500,695.10 |
| 2008-04-22 | THE ARCOP GROUP / GERSOVITZ MOSS | 423 - Engineering Consultants - Other | $67,145,777.15 |
| 2008-04-22 | CAPITAL ELEVATOR LTD. | 859 - Other Business Services not Elsewhere Specified | $73,651.02 |
| 2008-09-04 | NORR LIMITED | 421 - Architectural Services | $40,100,544.85 |
| 2009-02-28 | R&R AUTOMATION INC. | 859 - Other Business Services not Elsewhere Specified | $1,077,864.27 |
| 2009-03-25 | CAPITAL ELEVATOR LTD. | 859 - Other Business Services not Elsewhere | $108,703.30 |

Inspector   Console   Debugger   {} Style Editor   Performance   Memory   Network   Storage

All   HTML   CSS   JS   XHR   Fonts   Images   Media   Flash   WS   Other                     Filter URLs

Within the panel there are a number of options are available.

Clicking on the small trash can icon will remove all request and response data from the panel.

The remaining options filter the information displayed. For example, clicking on HTML limits the display to markup files such as HTML and XML. The CSS option limits the display to CSS files, the JS option to script files, the Images option to image files, and so on. The XHR option will show AJAX requests and responses.

The All option shows all request and response activity.

We'll look more closely at the HTML and XHR options.

**Watching network traffic**

The network panel is populated as soon as you take an action, or a script triggers an action, that makes a request to an external server. For example, if we click on the details link for IGF Vigilance Inc., seen above, the panel will display all information about requests and responses that follow the click.

**Contract Details**

PWGSC Services
Transparency
Proactive Disclosure
Disclosure of Contracts
Disclosure of Travel and Hospitality Expenses
Disclosure of Annual Expenditures for Travel, Hospitality and Conferences
Disclosure of Position Reclassifications
Disclosure of Grants and Contributions
Disclosure of wrongdoings in the workplace

N.B.: The contract date represents the date that the contract is recorded in the departmental financial system.

Details for a specific contract

| | |
|---|---|
| Vendor Name | IGF VIGILANCE INC. |
| Reference Number | 700006186 |
| Contract Date | 2008-04-22 |
| Description of Work | 420 - Engineering Services not Elsewhere Specified |

Inspector | Console | Debugger | {} Style Editor | Performance | Memory | Network | Storage

All HTML CSS JS XHR Fonts Images Media Flash WS Other — Filter URLs

| Status | Method | File | Domain | Cause | Type | Transferred | Size |
|---|---|---|---|---|---|---|---|
| 200 | GET | cl.pl?lang=eng;SCR=D;Sort=0;PF=CL201516Q... | www.tpsgc-pwgsc.gc.ca | document | html | 16.64 KB | 16.64 KB |
| 200 | GET | jquery.min.js | www.tpsgc-pwgsc.gc.ca | script | x-j | 91.44 KB | 91.44 KB |
| 200 | GET | util-min.css | www.tpsgc-pwgsc.gc.ca | stylesheet | css | 85.91 KB | 85.91 KB |
| 200 | GET | pe-ap-min.css | www.tpsgc-pwgsc.gc.ca | stylesheet | css | 85.65 KB | 85.65 KB |
| 200 | GET | theme-min.css | www.tpsgc-pwgsc.gc.ca | stylesheet | css | 93.15 KB | 93.15 KB |
| 200 | GET | sig-eng.gif | www.tpsgc-pwgsc.gc.ca | img | gif | 1.56 KB | 1.56 KB |
| 200 | GET | wmms.gif | www.tpsgc-pwgsc.gc.ca | img | gif | 3.09 KB | 3.09 KB |
| 200 | GET | piwik-min.js | www.tpsgc-pwgsc.gc.ca | script | x-j | 523 B | 523 B |
| 200 | GET | theme-min.js | www.tpsgc-pwgsc.gc.ca | script | x-j | 4.99 KB | 4.99 KB |
| 200 | GET | settings.js | www.tpsgc-pwgsc.gc.ca | script | x-j | 669 B | 669 B |
| 200 | GET | pe-ap-min.js | www.tpsgc-pwgsc.gc.ca | script | x-j | 99.53 KB | 99.53 KB |
| 200 | GET | piwik.js | stats.tpsgc-pwgsc.gc.ca | script | js | 17.02 KB | 43.08 KB |
| 200 | GET | nav_mega-mega_nav-eng.inc | www.tpsgc-pwgsc.gc.ca | xhr | plain | 5.85 KB | 5.85 KB |

21 requests | 724.34 KB / 519.21 KB transferred | Finish: 1.70 s | DOMContentLoaded: 860 ms | load: 1.71 s

As you can see, the browser has gone to the page and the panel below has been populated with information about 21 requests made. Below, there is a summary of the requests which shows that the total size of the request and response traffic was about 724 KB. Of that, about 519 was transferred from the remote server. When you visit a website, your browser will cache some of the response files on your hard drive so that when you return to the site, the files can be loaded quickly without having to request them from the server again.

The File column shows the URL to which the request was made. It's a relative URL, but if you hover over it you will see the full URL. The Method column shows the request method, GET or POST.

If you <Ctrl> Click (right click on a PC) on the file cell for a request, a context sensitive menu appears that allows you to copy the URL to paste into a browser's address bar.

Inspector | Console | Debugger | {} Style Editor | Performance | Memory | Network | Storage

All HTML CSS JS XHR Fonts Images Media Flash WS Other

Copy ▶ | Copy URL
Save All As HAR | Copy URL Parameters
| Copy as cURL
Edit and Resend | 
| Copy Request Headers
Open in New Tab | Copy Response Headers
Start Performance Analysis... | Copy Response
| Copy All As HAR

21 requests | 724.34 KB / 519.21 KB transferred | Finish: 1.70 s | DOMContentLoaded: 860 ms | load: 1.71 s

This feature is useful for examining files that are included as part of a response, such as JavaScript files, CSS style files, images, and other files.

The Status column reports on the status of the HTTP request. The code 200 means the request was successful. 304 indicates that the file was fetched from the browser cache. 404 means the resource did not exist on the server. You can read about all of the different possible status codes here: https://www.w3.org/Protocols/rfc2616/rfc2616-sec10.html

The Domain column indicates which Internet domain responded to the request.

The Size column indicates the size of the request and the returned response, in Bytes.

The time scale on the far right tells you how long the request took in milliseconds.

Clicking on one of the request rows will open a request details panel to the right (it can also be opened with the ⬛ icon. It has a number of tabs of its own. There's a lot of technical stuff here, but we'll make note of a few things.



Under the Headers tab, the Remote Address entry indicates the IP address of the server that sent the response. This can also be copied and pasted into the address bar of a web browser. We can also see the HTML request and response headers. HTML headers contain information used by the client and server computers. They are hidden from users in the normal day-to-day use of a browser.

Let's look at the request headers first.

The Host is the domain of the server. The referrer is the page from which the request came (we clicked on a link). Most important for our purposes is the User Agent. This is a character string that represents the type of browser, including the version, that made the request. Some web servers will

ignore requests that don't come from actual browsers. It is possible to set the header string sent by a script so it will appear to be a real browser. You can copy the user agent string from this panel, and paste it into a script. Of course, doing so raises potential ethical issues, and these are discussed in Chapter 9.

The response headers tell us the type of content in the response, in this case an HTML page, the character encoding, in this case ISO-859-1, the time and date of the request given in coordinated universal time, or GMT, and other technical details.

The Params tab contains the parameters that were passed to the external web server as part of the URL for the request. This can be extremely useful for understanding what you might need to send to the server in a scraping script.

| Headers | Cookies | Params | Response | Timings | Stack Trace | Security |
|---------|---------|--------|----------|---------|-------------|----------|

▽ Filter request parameters

▼ Query string

    lang: eng;SCR

The Response tab contains the actual content of the response, in this case the HTML page that was sent back.

| Headers | Cookies | Params | Response | Timings | Stack Trace | Security |
|---------|---------|--------|----------|---------|-------------|----------|

▼ Response payload

```
 1    <!DOCTYPE html>
 2 ▼  <!--[if IE 7]><html lang="en" class="no-js ie7"><![endif]-->
 3 ▼  <!--[if IE 8]><html lang="en" class="no-js ie8"><![endif]-->
 4    <!--[if gt IE 8 | !(IE)]><!-->
 5 ▼  <html lang="en" class="no-js">
 6 ▼  <!--<![endif]-->
 7 ▼  <head>
 8    <meta charset="utf-8" />
 9 ▼  <!-- Web Experience Toolkit (WET) / Bote  outils de l'exprience Web (BOEW)
10    wet-boew.github.com/wet-boew/License-eng.txt / wet-boew.github.com/wet-boew/Licenc
11 ▼  <!-- WET 3.0, PWGSC 1.0 file: 2col-eng.html -->
12 ▼  <!-- TitleStart -->
13 ▼  <!-- InstanceBeginEditable name="doctitle" -->
14 ▼  <title>Contract Details - Disclosure of Contracts - Proactive Disclosure - PWGSC</
15 ▼  <!-- InstanceEndEditable -->
16 ▼  <!-- TitleEnd -->
17 ▼  <!-- MetadataStart -->
18 ▼  <!-- InstanceBeginEditable name="metadata" -->
19
20     <meta name="dcterms.title" content="Contract Details - Disclosure of Contracts -
21        <meta name="dcterms.subject" content="contracts; service contracts; federal co
22        <meta name="dcterms.language" title="ISO639-2" content="eng" />
23        <meta name="dcterms.creator" content="Government of Canada, Public Works and G
24        <meta name="dcterms.publisher" content="Public Works and Government Services C
25        <meta name="dcterms.issued" title="W3CDTF" content="2004-10-01" />
26        <meta name="dcterms.modified" title="W3CDTF" content="2017-08-15" />
27        <meta name="dcterms.description" content="Proactive disclosure of contracts" /
28     <meta name="description" content="Proactive disclosure of contracts" />
29     <meta name="keywords" content="proactive disclosure, contract, reports, disclosure
30
31 ▼  <!-- InstanceEndEditable -->
```

This may seem superfluous when the response is a plain HTML page, but sometimes you will want to see the response made to an XHR (XMLHttpRequest) request when AJAX is being used, the data sent back in response to a search request, or the contents of a JavaScript or CSS file.

The Timings tab provides information on how long the request and response took.

**When a page isn't completely reloaded**

As discussed in Chapter 9, not all web sites request a whole new page each time a user interacts with the page. Instead, data may be requested and passed back to the browser without having to reload the page. This can be done in different ways. For example, a request may be made for an external html resource that is then placed in the main page by way of an iframe. Or a site may use a technology called AJAX. When these kinds of methods are being used, opening and looking at the HTML source for a page will only show you the original HTML that was sent by the server. Any changes to the DOM (document object model), which could be thought of as the current state of the web page the browser has stored in memory, will not be reflected. The network panel, however, can allow you to peek under the hood, and see everything that is happening.

Let's say you wanted to have a look at flight arrivals and departures and John F. Kennedy Airport in New York. You could go to the website maintained by the Port Authority of New York and New Jersey at https://www.panynj.gov/airports/flight-status.html

You may notice that the actual flight status information, which comes from an outside provider, is slightly delayed in loading compared to the rest of the page. This is because you are being provided with the latest information, and it is being added after the main page loads. Here is what the fully loaded page, with the embedded flight information looks like:

If we use the element inspector to examine the area of the page that contains the flight information, we can see that contained within the <div> tag with the id of "viewport" there is an iframe that has as its source, a different URL. It's a relative URL, so we can't see the domain, but we can see that by using the network monitor.
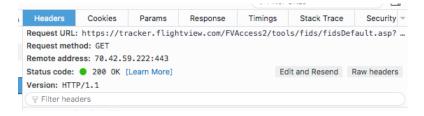


We'll switch to the network tab in developer tools, then choose the HTML filter (if you go straight to the network tab after opening developer tools, you may need to reload the page in the browser before the network panel populates). We're using the HTML filter because an iframe is used to embed an HTML document.



We can see that four files were transferred, including one that is almost a megabyte in size. We'll guess that might be the file we want, because there's quite a bit of arrivals data.

If we now open the request details panel by clicking on the row for the request, we can see that the request was made to a URL that does not appear in the address bar of our web browser.

This is the full URL:

https://tracker.flightview.com/FVAccess2/tools/fids/fidsDefault.asp?accCustId=PANYNJ&fidsId=20001&fidsInit=arrivals&fidsApt=JFK

Now, if we click on the Response tab, we can see the HTML filed that was transferred.

```
77    </td>
78    <td>5011</td>
79    <td>Boston, MA (BOS)</td>
80    <td>Arrived</td>
81    <td><!--dtDateTime(2017-08-16,16:23:00,TwelveHour)-->4:23 PM</td>
82    <td><!--dtDateTime(2017-08-16,16:44:00,TwelveHour)-->4:44 PM</td>
83    <td>Term 5 - 03</td>
84    <td><a href="javascript:void(0)" onclick="ffDtNm('fA','5011','EI','BOS','
85    </tr>
86    <tr class="odd">
87    <td class="c1">
88    <table class="ffAlTbl">
89    <tr>
90    <td id="ffAlLog" class="ffAlLog"><script>ffWrAlLg("EI");</script></td>
91    <td id="ffAlLbl" class="ffAlLbl">Aer Lingus</td>
92    </tr>
93    </table>
94    </td>
95    <td>5013</td>
96    <td>Boston, MA (BOS)</td>
97    <td>Scheduled</td>
98    <td><!--dtDateTime(2017-08-16,20:32:00,TwelveHour)-->8:32 PM</td>
99    <td><!--dtDateTime(2017-08-16,20:25:00,TwelveHour)-->8:25 PM</td>
100   <td>Term 5 - 01</td>
101   <td><a href="javascript:void(0)" onclick="ffDtNm('fA','5013','EI','BOS','
102   </tr>
103   <tr class="even">
104   <td class="c1">
105   <table class="ffAlTbl">
106   <tr>
107   <td id="ffAlLog" class="ffAlLog"><script>ffWrAlLg("EI");</script></td>
108   <td id="ffAlLbl" class="ffAlLbl">Aer Lingus</td>
109   </tr>
```

If we like, we can copy and paste the HTML into a separate file, if you wished to retain it for future reference.
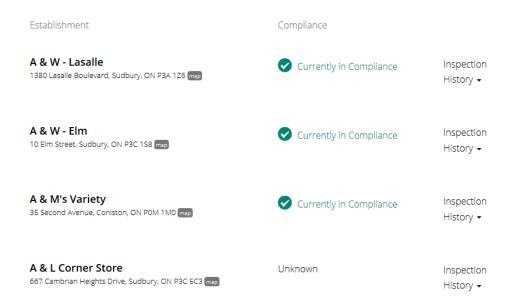
## Looking at AJAX

As we discussed above, another way used to insert information into an existing web page is the AJAX protocol.

We'll take a look at a site used by the Sudbury District Health Unit in the northern part of the Canadian province of Ontario to provide information on restaurant inspections to the public. The page now contains archived information, and will be replaced by a new site, according to the municipality.
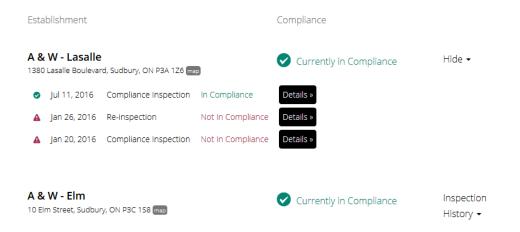
The site showed the restaurants in the city, and whether each is in compliance. You could also access a more detailed inspection history for each premise.

## 1743 Establishments Found

| Establishment | Compliance | |
|---|---|---|
| **A & W - Lasalle**<br>1380 Lasalle Boulevard, Sudbury, ON P3A 1Z6 `map` | ✅ Currently in Compliance | Inspection History ▾ |
| **A & W - Elm**<br>10 Elm Street, Sudbury, ON P3C 1S8 `map` | ✅ Currently in Compliance | Inspection History ▾ |
| **A & M's Variety**<br>35 Second Avenue, Coniston, ON P0M 1M0 `map` | ✅ Currently in Compliance | Inspection History ▾ |
| **A & L Corner Store**<br>667 Cambrian Heights Drive, Sudbury, ON P3C 5C3 `map` | Unknown | Inspection History ▾ |

If you click on the Inspection History link for a premise, the page was populated with individual inspections for that facility, complete with more links to delve down to more precise detail.

## 1743 Establishments Found

| Establishment | | | Compliance | |
|---|---|---|---|---|
| **A & W - Lasalle**<br>1380 Lasalle Boulevard, Sudbury, ON P3A 1Z6 `map` | | | ✅ Currently in Compliance | Hide ▾ |
| ✅ | Jul 11, 2016 | Compliance Inspection | In Compliance | Details » |
| ⚠️ | Jan 26, 2016 | Re-inspection | Not In Compliance | Details » |
| ⚠️ | Jan 20, 2016 | Compliance Inspection | Not In Compliance | Details » |
| **A & W - Elm**<br>10 Elm Street, Sudbury, ON P3C 1S8 `map` | | | ✅ Currently in Compliance | Inspection History ▾ |

An examination of the page source, however, shows that the detail listed is not present. The DOM of the page is being manipulated using JavaScript that runs based on user interactions, such as clicking on an Inspection History link. New data is being added to the page dynamically.

We can have a closer look at what is going on using the Network monitor in Firefox's developer tools.

When we first load the page, we can see that five requests were made to the server, totalling 529 kb of data sent and received.
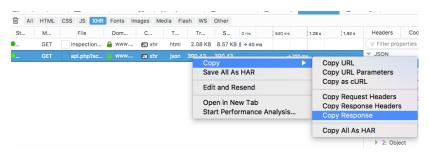


The site used to serve up a great many more resources, but has been reduced in scope as it only archives date up till January 2017.

As you can see, two of the files transferred have XHR under type. XHR request is another way of saying AJAX request.

If we switch to the XHR tab, we can see just these two and we can also see that one of the files is an html file and other a JSON file. JSON stands for JavaScript Object Notation, which as discussed in Chapter 2 is a plain-text data format often used when the data is intended to be machine read.

If we want to see the whole response, we can <ctrl> click (right click on a Windows computer) on the row for the JSON file in the network panel, and choose Copy>Response.



We can then paste the resulting JSON into a separate file to examine further. You can then use a Python script or an online converter to convert the JSON to a CSV file for import into a spreadsheet or database program. In this way, you are using the web development tool as a rudimentary web scraping tool. For a few pages, it could be all you need.

If we like, we can also look at the response in the request details window. This allows us to examine each individual object in the JSON file.

An object in this context is a data construct that holds all of the information for one row of data, that is, one facility that was inspected.

If we click on the expansion triangle for object 0, we can see the data contained in that object, neatly organized.



This makes it trivially easy to examine the data structure.

There is a great deal more that the Firefox developer tools, and those in other browsers, can do, including helping you develop and debug JavaScript, manipulate the HTML and CSS of the page to see, in real time, how the change would affect the page, and much more that we simply couldn't cover here. These are remarkable tools and belong in every data journalist's toolkit.