

# STAT 209

## Scraping Data from HTML Tables

March 28, 2018

Colin Reimer Dawson

# Data from Web Pages

Sometimes you see data online not in the form of a file:

## THANKSGIVING WEEKENDS (1982-Present)

3-day Openings

3-day All Movies

5-day Openings

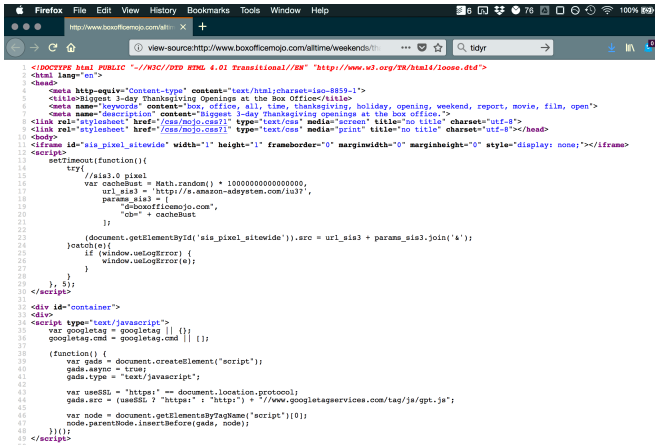
5-day All Movies

Rank	Title (click to view)	Studio	Opening*	% of Total	Theaters	Avg.	Total Gross^	Release Date%
1	<b>Frozen</b>	BV	<b>\$67,391,326</b>	16.8%	3,742	\$18,009	\$400,738,009	11/27/13
2	<b>Toy Story 2</b>	BV	<b>\$57,388,839</b>	23.3%	3,236	\$17,734	\$245,852,179	11/24/99
3	<b>Moana</b>	BV	<b>\$56,631,401</b>	22.8%	3,875	\$14,615	\$248,757,044	11/23/16
4	<b>Coco</b>	BV	<b>\$50,802,605</b>	24.3%	3,987	\$12,742	\$209,372,919	11/22/17
5	<b>Tangled</b>	BV	<b>\$48,767,052</b>	24.3%	3,603	\$13,535	\$200,821,936	11/24/10
6	<b>The Good Dinosaur</b>	BV	<b>\$39,155,217</b>	31.8%	3,749	\$10,444	\$123,087,120	11/25/15
7	<b>Enchanted</b>	BV	<b>\$34,440,317</b>	26.9%	3,730	\$9,233	\$127,807,262	11/21/07
8	<b>101 Dalmatians (1996)</b>	BV	<b>\$33,504,025</b>	24.6%	2,794	\$11,991	\$136,189,294	11/27/96
9	<b>A Bug's Life</b>	BV	<b>\$33,258,052</b>	20.4%	2,686	\$12,382	\$162,798,565	11/25/98
10	<b>Four Christmases</b>	WB (NL)	<b>\$31,069,826</b>	25.9%	3,310	\$9,387	\$120,146,040	11/26/08

Figure: Source: Box Office Mojo ([link](#))

# Messy Embedding

The data is embedded (somewhere) in the HTML code of the site



```
1 <!DOCTYPE html PUBLIC "-//W3C//DTD XHTML 1.0 Transitional//EN" "http://www.w3.org/TR/xhtml1/loose.dtd">
2 <html lang="en">
3 <head>
4 <meta http-equiv="Content-type" content="text/html; charset=iso-8859-1">
5 <title>Biggest 3-day Thanksgiving Openings at the Box Office</title>
6 <meta name="keywords" content="box, office, all, time, thanksgiving, holiday, opening, weekend, report, movie, film, open">
7 <meta name="description" content="Biggest 3-day Thanksgiving openings at the box office.">
8 <link rel="stylesheet" href="/css/mojo.css?" type="text/css" media="screen" title="no title" charset="utf-8">
9 <link rel="stylesheet" href="/css/mojo.css?" type="text/css" media="print" title="no title" charset="utf-8"></head>
10 <body>
11 <iframe id="sis_pixel_sitewide" width="1" height="1" frameborder="0" marginwidth="0" marginheight="0" style="display: none;"></iframe>
12 <script>
13 setTimeout(function(){
14     try{
15         //sis3.0 pixel
16         var cacheBust = Math.random() * 10000000000000000;
17         url_sis3 = 'http://s.smaron-adsystem.com/ia3?';
18         params_sis3 = {
19             "d":"boxoffice/mojo.cca",
20             "cb" + cacheBust
21         };
22     }
23     (document.getElementById('sis_pixel_sitewide')).src = url_sis3 + params_sis3.join('&');
24 }catch(e){
25     if (window.useLogError) {
26         window.useLogError(e);
27     }
28 }, 5);
29 </script>
30 <div id="container">
31 <div>
32 <script type="text/javascript">
33 var googletag = googletag || {};
34 googletag.cmd = googletag.cmd || [];
35 (function() {
36     var gads = document.createElement("script");
37     gads.async = true;
38     gads.type = "text/javascript";
39     var useSSL = "https:" == document.location.protocol;
40     gads.src = (useSSL ? "https:" : "http:") + "//www.googletagservices.com/tag/js/gpt.js";
41     var node = document.getElementsByTagName("script")[0];
42     node.parentNode.insertBefore(gads, node);
43 })();
44 </script>
```

# Our Goals

Get the data from the web into R in a *reproducible* way. That means:

1. No manual editing of text
2. No copying and pasting
3. No fiddling with the data in Excel

In other words, we need to go from what's on the web to a data frame (and visualizations) using a *self-contained* script

## More tidyverse Building Blocks

