

# STAT 209

## Scraping and Cleaning Data from the Web

November 5, 2019

Colin Reimer Dawson

# Data from Web Pages

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

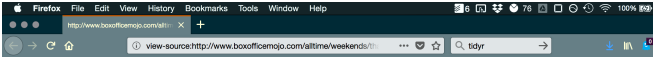
## THANKSGIVING WEEKENDS (1982-Present)

	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 HTML 4.01 Transitional//EN" "http://www.w3.org/TR/html4/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_sideview" 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.amazon-adsystem.com/uis3",
18       params_sis3 = {
19         "d":"boxoffice Mojo.com",
20         "cb":" + cacheBust
21       };
22
23       (document.getElementById('sis_pixel_sideview')).src = url_sis3 + params_sis3.join('&');
24     }catch(e){
25       if (window.useLogError) {
26         window.useLogError(e);
27       }
28     }
29   }, 5);
30 </script>
31
32 <div id="container">
33 <div>
34 <script type="text/javascript">
35   var googletag = googletag || {};
36   googletag.cmd = googletag.cmd || [];
37
38   (function() {
39     var gads = document.createElement("script");
40     gads.async = true;
41     gads.type = "text/javascript";
42
43     var useSSL = "https:" == document.location.protocol;
44     gads.src = (useSSL ? "https:" : "http:") + "//www.googletagservices.com/tag/js/gpt.js";
45
46     var node = document.getElementsByTagName("script")[0];
47     node.parentNode.insertBefore(gads, node);
48   })();
49 </script>
50
```

# Our Goals

Get the data from the web into our visualizations 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



## Necessary Steps

1. Pulling the HTML from the site
2. Extracting the tables
3. Converting a table to a data frame
4. Fixing variable formats (quantitative data is numeric, dates are parsed)
5. Visualization

We also want to be able to do some simple text substitutions to clean up the ways variables are encoded; particularly if we are counting occurrences of certain values.

# Lab

- Today's lab just scratches the surface of web scraping and data cleaning.
- We'll revisit some of these topics in more detail later (time-permitting)