# Web scraping in R using the rvest package

R Ladies Meetup
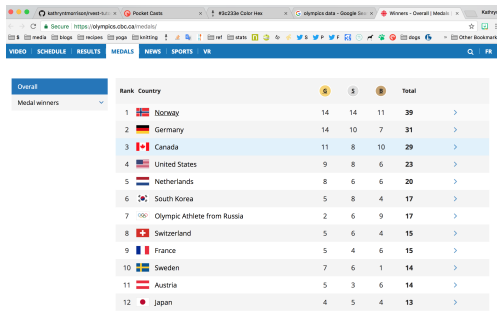
Kathryn Morrison, Co-Founder
Precision Analytics
March 15, 2018

## What is web scraping and why is it useful?

- Say you want to download some data off a website to do some analysis:



- You could manually copy and paste the data into a spreadsheet, but it would take a long time.
- Web scraping is a faster alternative to this

**What is web scraping and why is it useful?**

- Web scraping is basically automating the collection/download of unstructured data from online sources into a format you can use, e.g., plain text, csv, JSON, etc
- You can use many languages (e.g., python, java, *R*!)

Some popular approaches to web scraping:

1. Text pattern matching with regular expressions
2. Using an API
3. Parsing a website (what we'll do)

- I'll go through a simple parsing example
- If you're using an API or something more complex, you may need more technical skills/tools than rvest

## Scraping in rvest

- Knowing some basics of HTML and CSS is helpful (but not required for the basics)
- You can use the SelectorGadget chrome extension to find the necessary tags

The rvest package was inspired by 'Beautiful Soup' in python

Rvest was created by Hadley Wickham (of course)

We'll go through a really simple example