

Mind the Gap: Leveraging Web Technologies from R

Carson Sievert

2015-11-03

Contents

| | | |
|----------|--------------------------|----------|
| 1 | Problem Statement | 1 |
| 2 | Overview | 1 |
| 3 | Scope | 4 |
| 4 | Background | 5 |
| | References | 5 |

1 Problem Statement

Web technologies offer exciting new opportunities to advance the field of statistical computing; particularly in the areas of data acquisition and statistical interactive graphics. In theory, freely available data on the web is accessible; but in practice, one needs an extensive knowledge of web technologies to acquire it. With better software for curating data on the web, we can reduce barriers to access, and enable a reproducible workflow.

Reproducibility and portability have remained a challenge for interactive statistical graphics for decades. Interactive statistical graphics software is often difficult to install since it typically assumes non-standard software is available on users' systems. Modern web browsers with **HTML5** support are now ubiquitous, and provide a rich ecosystem for interactive graphics. However, interfacing this ecosystem with statistical software remains difficult and, depending on the application, can require a myriad technologies. Again, with better software, we can reduce the startup costs involved with producing web-based interactive statistical graphics.

The current trend in web-based interactive statistical graphics is provide various language bindings to **JavaScript** charting libraries. To test whether the entire software stack is working as intended, it's common to verify properties of the data sent to the binding, but this does not guarantee that the end result is what we expect. A proper testing framework for this type of software should be able to construct and manipulate the Document Object Model (DOM) using technologies available to modern web browsers. To our knowledge, **animint** is the first R package to implement this testing approach, and some of the lessons learned could be used to construct a more reliable and easier to use testing suite.

2 Overview

This proposal is a collection of work towards making web-based interactive statistical graphics and data on the web more accessible to R programmers. I currently maintain packages 7 R packages in this direction (and have co-authored a handful of others): **pitchRx**, **bbsscraperR**, **XML2R**, **rdom**, **LDavis**, **animint**, **plotly**. This section provides a high-level overview of this work within R's broader ecosystem and points out holes remaining to be filled. It also describes portions that I've contributed to in joint work such as **LDavis**, **animint** and **plotly**.

The R packages **pitchRx**, **bb scrapeR**, **XML2R**, and **rdom** all provide utilities for acquiring data hosted on the web. Amongst them, **pitchRx** and **bb scrapeR** have the highest level interface and are aimed at a known set of Uniform Resource Locators (URLs). Their interface is designed so that users do not need any understanding of underlying file formats that contain the data (e.g., **XML**, **JSON**). In order to acquire some data, users just need the package installed, so a very minimal amount of computational setup is required. If acquiring large amounts of data, to avoid memory limitations, users may divert incoming data to a database using any valid R database connection (Databases 2014).

For these reasons, **pitchRx** and **bb scrapeR** provide a nice resource for teaching and practicing applied statistics, and serve as a model for providing access to clean versions of messy datasets on the web (Unwin 2010). Providing *access* to data in this way is more desirable than rehosting data for several reasons. In some cases, it helps avoid legal issues with rehosting copyrighted data. Furthermore, the packages are self-documenting, so users can inspect the cleaning and transformations performed on the data to ensure its integrity. They are also versioned which makes the acquisition, and thus any downstream analysis, more reproducible and transparent. Perhaps most importantly, the data that these packages provide access to are updated at least several times a day, so users can keep their local copies up to date without any assistance from the maintainer.

pitchRx and **bb scrapeR** are specific to sports data, but the same model is popular across many different domains. Perhaps the largest centralized effort in this direction is lead by **rOpenSci**, a community of R developers that, at the time of writing, maintains more than 50 packages providing access to scientific data ranging from bird sightings, species occurrence, and even text/metadata from academic publications. This provides a tremendous service to researchers who want to spend their time building models and deriving insights from data, rather than learning the programming skills necessary to acquire and clean it.

It's becoming increasingly clear that "meta" packages that standardize the interface to data acquisition/curation in a particular domain would be tremendously useful. However, it is not clear how such interfaces should be designed. The **etl** package (a joint work with Ben Baumer) is one step in this direction and actually aims to provide a standardized interface for *any* data access package that fits into an Extract-Transform-Load paradigm (Baumer and Sievert). The package provides generic **extract/transform/load** functions, but requires developers to write **extract/transform** methods for the specific data source. In theory, the default **load** method works for any application; as well as database management operations such as **update** and **clean**.

The Web Technologies and Service CRAN Task View does a great job of summarizing data access packages in R, and even breaks them down by their domain application (e.g., Government, Finance, Earth Science, etc) (Scott Chamberlain). It also details more general tools for requesting, parsing, and working with popular file formats on the web from R that many of the data access packages use to implement their functionality. Two other packages I maintain: **XML2R** and **rdom** fit into this broader category.

pitchRx and **bb scrapeR** leverage **XML2R**: a wrapper around the **XML** package for transforming XML content into tables (Lang and CRAN Team 2015). **XML** provides low-level R bindings to the libxml2 C library for parsing XML content (Veillard 2006). **XML2R** builds on this functionality and makes it possible to acquire, filter, and transform XML content into table(s) without any knowledge of the verbose **XPath** and **XSLT** languages. These high-level semantics make it easier to maintain projects such as **pitchRx** and **bb scrapeR** since it drastically reduces the amount of code. Chapter Taming PITCH/x Data with XML2R and **pitchRx** details these abstractions, and how they are used in **pitchRx**, in more detail.

rdom makes it easy to render dynamic web pages and access the Document Object Model (DOM) from R via the headless browsing engine phantomjs (Sievert 2015). This fills a void where other web scraping packages in R (e.g., **XML**, **xml2**, **rvest**) currently fall short. These packages can download, parse, and extract bits of the HTML page source, which is static, but they lack a browsing engine to fully render the DOM. If the DOM cannot be rendered, content that is dynamically generated (e.g., with **JavaScript**) cannot be acquired. The R package **RSelenium** can also render dynamic web pages and simulate user actions, but its broad scope and heavy software requirements make it harder to use and less reliable compared to **rdom**.

The **animint** package uses **RSelenium** in order to test whether visualizations render correctly in a web browser. When I started on **animint**, we were only testing the R code that compiles **ggplot2** objects as

JSON objects, but had no way to programmatically test the **JavaScript** code that takes the JSON as input. If **animint** were limited to static web-based graphics, we could use the lighter-weight **rdom** for testing, but we need to simulate user actions to verify animations and interactive features work as expected.

In addition to adding infrastructure for testing **animint**'s renderer, I've made a number of other contributions:

1. Wrote bindings for embedding **animint** plots inside of knitr/rmarkdown/shiny documents, before the advent of **htmlwidgets**, which provides standard conventions for writing such bindings (Vaidyanathan et al. 2015). At the time of writing, **htmlwidgets** can only be rendered from the R console, the R Studio viewer, and using R Markdown (v2). For this reason, we decide to not use **htmlwidgets** since users may want to incorporate this work into a different workflow.
2. Wrote **animint2gist**, which uploads an **animint** visualization as a GitHub gist, which allows users to easily share the visualizations with others via a URL link.
3. Implemented **ggplot2** facets (i.e., **facet_wrap** and **facet_grid**) as well as the fixed coordinate system (i.e., **coord_fixed**).
4. Mentored and assisted Kevin Ferris during his 2015 Google Summer of Code project where he implemented theming options (i.e., **theme**), legend interactivity, and selectize widgets for selecting values via a drop-down menu.

When I started on **plotly**, its core functionality and philosophy was very similar to **animint**: create interactive web-based visualizations using **ggplot2** syntax (Sievert et al.). However, **plotly**'s **JavaScript** graphing library supports chart types and certain customizations that **ggplot2**'s syntax doesn't support. Realizing this, I initiated and designed a new domain-specific language (DSL) for using **plotly**'s **JavaScript** graphing library from R. Although its design is inspired by **ggplot2**'s **qplot** syntax, the DSL does not rely on **ggplot2**, which is desirable since its functionality won't break when **ggplot2** internals change.

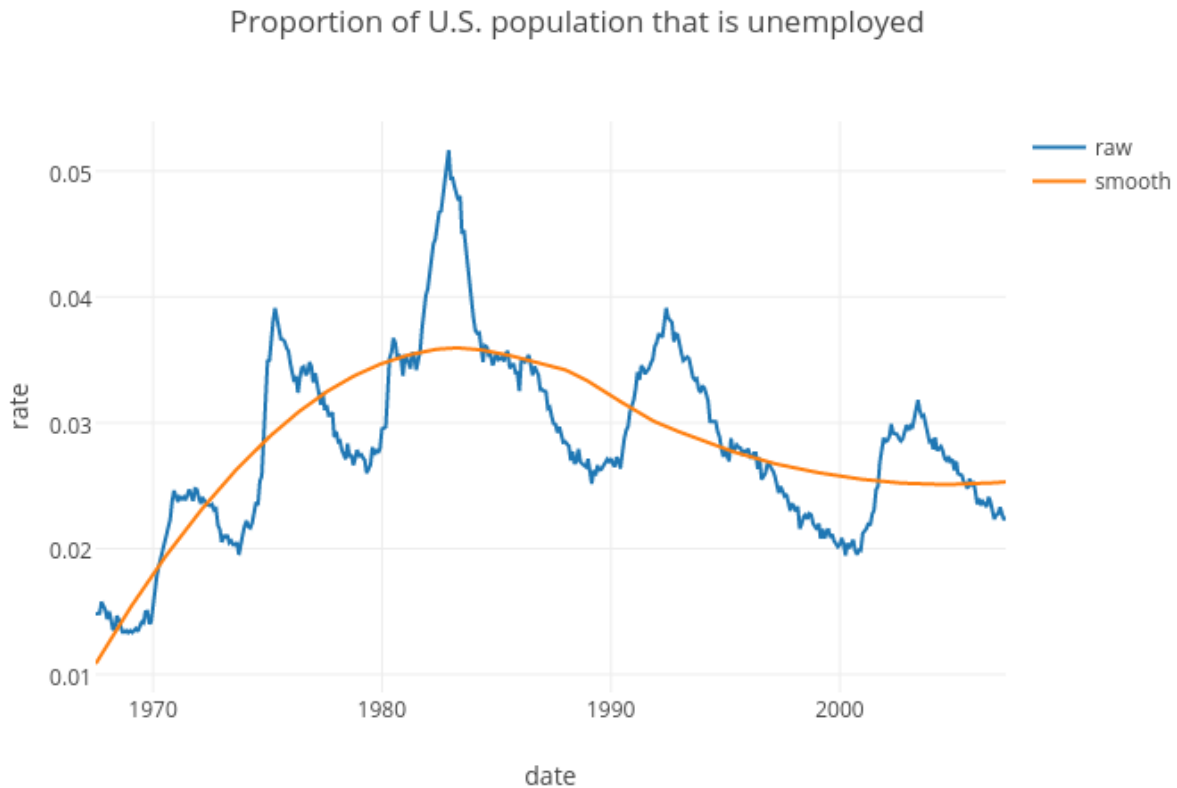
plotly's 'native' R DSL is heavily influenced by concepts deriving from pure functional programming. The output of a pure function is completely determined by its input(s), and because we don't need any other context about the state of the program, it's easy to read and understand the intention of any pure function. When a suite of pure functions are designed around a central object type, we can combine these simple pieces into a pipeline to solve complicated tasks, as is done in many popular R packages such as **dplyr**, **tidyr**, **rvest**, etc (Wickham 2015).

plotly's pure functions are deliberately designed around data frames so we can conceptualize a visualization as a pipeable sequence of data transformations, model specifications, and mappings from the data/model space to visual space. With the R package **ggvis** (Chang and Wickham 2015), one can also mix data transformation and visual specifications in a single pipeline, but it does so by providing S3 methods for **dplyr**'s generic functions, so all data transformations in a **ggvis** pipeline have to use these generics. By directly modeling visualizations as data frames, **plotly** removes this restriction that transformation must derive from a generic function, and removes the burden of exporting transformation methods on its developers.

plotly even respects transformations that remove attributes used to track visual properties and data mappings. To demonstrate, in the example below, we plot the raw time series with **plot_ly()**, fit a local polynomial regression with **loess()**, obtain the observation-level characteristics of the model with **augment()** from the **broom** package, layer on the fitted values to the original plot with **add_trace()**, and add a title with **layout()**.

```
library(plotly)
library(broom)
economics %>%
  transform(rate = unemploy / pop) %>%
  plot_ly(x = date, y = rate, name = "raw") %>%
  loess(rate ~ as.numeric(date), data = .) %>%
```

```
augment() %>%
add_trace(y = .fitted, name = "smooth") %>%
layout(title = "Proportion of U.S. population that is unemployed")
```



To make this possible, a special environment within **plotly**'s namespace tracks not only visual mappings/properties, but also the order in which they are specified. So, if a **plotly** function used to modify a visualization (e.g., `add_trace()` or `layout()`) receives a data frame without any special attributes, it retrieves the last plot created, and modifies that plot.

animint and **plotly** could be classified as general purpose software for web-based interactive and dynamic statistical graphics; whereas **LDavis**, could be classified as software for solving a domain specific problem. The **LDavis** package creates an interactive web-based visualization of a topic model fit to a corpus of text data using Latent Dirichlet Allocation (LDA) to assist in interpretation of topics. The visualization itself is written entirely with HTML5 technologies and makes use of the JavaScript library d3js (Heer 2011) to implement advanced interaction techniques that higher-level tools such as **plotly**, **animint**, and/or **shiny** do not currently support.

3 Scope

This section describes work to be achieved before completion of the thesis.

4 Background

Go into more detail?

References

- Baumer, Ben, and Carson Sievert. *Etl: Extract-Transfer-Load Framework for Medium Data*. <http://github.com/beanumber/etl>.
- Chang, Winston, and Hadley Wickham. 2015. *Ggvis: Interactive Grammar of Graphics*. <http://CRAN.R-project.org/package=ggvis>.
- Databases, R Special Interest Group on. 2014. *DBI: R Database Interface*. <http://CRAN.R-project.org/package=DBI>.
- Heer, Michael Bostock AND Vadim Ogievetsky AND Jeffrey. 2011. “D3: Data-Driven Documents.” *IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis)*. <http://vis.stanford.edu/papers/d3>.
- Lang, Duncan Temple, and the CRAN Team. 2015. *XML: Tools for Parsing and Generating XML Within R and S-Plus*. <http://CRAN.R-project.org/package=XML>.
- Scott Chamberlain, Patrick Mair, Thomas Leeper. “CRAN Task View: Web Technologies and Services.” <http://cran.r-project.org/web/views/WebTechnologies.html>.
- Sievert, Carson. 2015. *Rdom: Access the DOM of a Webpage as HTML Using Phantomjs*. <https://github.com/cpsievert/rdom>.
- Sievert, Carson, Chris Parmer, Toby Hocking, Scott Chamberlain, Karthik Ram, Marianne Corvellec, and Pedro Despouy. *Plotly: Create Interactive Web-Based Graphs via Plotly’s API*. <https://github.com/ropensci/plotly>.
- Unwin, Antony. 2010. “Datasets on the web: a resource for teaching Statistics?” *MSOR Connections*, November, 1–4.
- Vaidyanathan, Ramnath, Yihui Xie, JJ Allaire, Joe Cheng, and Kenton Russell. 2015. *Htmlwidgets: HTML Widgets for R*. <http://CRAN.R-project.org/package=htmlwidgets>.
- Veillard, Daniel. 2006. “Libxml: The XML c Parser and Toolkit of Gnome Parsing.” <http://www.xmlsoft.org>.
- Wickham, Hadley. 2015. “Pipelines for Data Analysis.” <http://bids.berkeley.edu/resources/videos/pipelines-data-analysis>.