# Data Collection from Syndicated Web Sources - APIs

Session 4 @ CBA Batch 12

April 2019

sudhir_voleti@isb.edu

ISB

1

---

## Session Outline

- Session 2 Continued: Webscraping Basics
  - [1] Atis.org revisited [2] Scraping Cricbuzz [3] Scraping Amazon Reviews

- Py Scraping Google search results

- Scraping Dynamic pages with Webdrivers
  - [1] Scraping python.org

- Intro to DC from APIs
  - DarkSky weather API, openMovieDB API

- DC from Finance APIs – Quandl

ISB

2

# Session 2 Webscraping Basics

# Continued

ISB

3

---

## Atis Glossary Scraping revisited

- **Background:**
- Atis.org website's changed in the past few months. Old codes weren't working.

- So I and my RA, Ms Anisha both rewrote the codes in parallel.

- LMS contained her version which was <u>different</u> from mine.

- Let me briefly walk you through my (simpler) version.

- Aim is to demo the use of user-defined functions in webscraping

- Pls open: **Webscraping basics with rvest.Rmd**

ISB

4

## Atis Glossary Scraping Recap

- First 2 steps were common to what we did previously.

- Step 3 onwards we see a change.

- We wrote 1 function for 1 letter and looped over all letter links.

- <u>More generally:</u>
- We write a functions (for 1 link) but that link contains a set of other (secondary) links.
- So we write another functions for those secondary links and so on.

- We'll use this approach going forward in the next 2 examples.

ISB

5

## Scraping IPL 2019 match schedules from Cricbuzz

- Open: **Another rvest exercise - cricbuzz.Rmd**

- The webpage looks like this - what all attributes can we easily scrape?

- Use SelectorGadget and ID the tags for those attributes. For example:

- See the photos tab on the page. How to scrape images?

- Exercise done. Learnings? Applications? Etc.

ISB

6

## Scraping Structured Webpages - Amazon reviews

- Why might anyone want to scrape product reviews from Amazon?

- Recall what structured webpages are.

- Consider how we discerned structure in Amazon review URLs.

- Then we wrote a function to scrape one page.

- We constructed URLs for required page range, and looped our func over it.

- Exercise done. Learnings? Applications? Etc.

ISB

7

# Py-scraping Google Search Results

- Suppose you're running a kebab joint... You want to know where you stand vis-à-vis local (city based) competition.

- Maybe you want to know what shows up when someone Google-searches for, say, 'Kebab Hyderabad' etc ...

- What pages would show up, in which order, what content would they contain, what sentiment would they express, etc.

- *Why* focus on Google search only? *

- And what all can you do once that data are in?*

- Python makes it super-simple to scrape Google search results → save it as a structured object → that you can analyze in R (for example).

ISB

8

Google Search Content for 'ISB CBA'

9



Extracting Google Search Results in Python

10

## Basic Text Analysis on the Scraped Text in R

```
19
20  data = read.csv(file.choose(), stringsAsFactors = F)
21  dim(data)
22  names(data)
23
24  # Remove PDF documents and links as they were not read correctly
25  if (length(grep('.pdf',data$url)) != 0) {
26      data = data[-grep('.pdf',data$url),]
27      dim(data)
28  }
29
30  data$text  =  iconv(data$text, "latin1", "ASCII", sub="")   # Keep only ASCII characters
31  wordCorpus <- Corpus(VectorSource(data$text))
32  wordCorpus <- tm_map(wordCorpus, removePunctuation)
33  wordCorpus <- tm_map(wordCorpus, content_transformer(tolower))
34  wordCorpus <- tm_map(wordCorpus, removeWords, stopwords("english"))
35  wordCorpus <- tm_map(wordCorpus, removeNumbers)
36  wordCorpus <- tm_map(wordCorpus, stripWhitespace)
37
38  pal <- brewer.pal(9,"YlGnBu")
39  pal <- pal[-(1:4)]
40  set.seed(123)
41  #windows()
42  wordcloud(words = wordCorpus, scale=c(4,0.5), max.words=100, random.order=FALSE,
43            rot.per=0.35, use.r.layout=FALSE, colors=pal)
44
45  ########################################################
```

ISB

11

# Webdrivers for Web-scraping

ISB

12

## Using Selenium in Py

- Open the HTML file 'Intro to webdrivers - Selenium in Py'

- Let's walk though it step by step.

- What I'll show next is fairly basic. However, ...

- If you're aware of alternatives, better ways to do the same thing etc., pls speak up and share with the class.

ISB

13

## Using Webdrivers: Recap

- What are webdrivers and where are they most used?

- What modules did we invoke for using webdrivers?

- What main functions were called? What did they do?

- What further possibilities come to mind with webdriver use?

- Ready for some basic homework involving py and selenium?

ISB

14

Some API

Preliminaries

ISB

15

## Some Common Secondary Data Sources

- Let's quickly organize common secondary data sources:

Internal sources

External sources

Records

Systems (e.g., ERP, CRM, SCM etc.)

Web

Public

So where do APIs fit in here?

Proprietary

Govt.

Website, apps etc.

Etc.

Syndicated Data Providers

ISB

16

## APIs: Some Preliminaries

- What is an API?

  Application Programming Interfaces are *interfaces* between 2 *services* …

- Examples?

- "There's an API for that".

- *Why* do firms like FB or Google put out APIs?

  Potential source of revenue since data is currency→ monetize the data asset; Invite developers to deploy cool stuff through their platform

- In which domains might APIs be likely to be found? How many might be there?

ISB

17

## APIs: The Growth Story across Domains



FASTEST GROWING API CATEGORIES SINCE 2014 - PRIMARY OR SECONDARY

ProgrammableWeb
GROWTH (SINCE 2014)

ISB

18

## API Preliminaries:  Data Storage Formats

- Unlike HTML & DOM which are more of data *markup and display* formats, JSON and XML are popular data *storage* formats.

- Consider an example of fields {Name, age, occupation} for two people A & B.
  - {Ravi, 38, Graphic Designer}
  - {Anu, Sales Executive, 27}

- Consider how a person vs  how a machine would read & understand.
  - Why the difference? What can be done about *ensuring* such doesn't happen?

- Enter data storage formats like JSON and XML.
  - These contain both the field names and the field values for every data point.
  - Verbose, but accurate.
  - Sample this example

ISB

19

## JSON & XML Data Storage Formats

- Here's a quick view of what JSON output looks like …

- Can you ID the field names (or 'keys') and values?

- And now a quick view of what XML output looks like …

- Can you ID the field names (or 'keys') and values?

- Note the ability to nest and build hierarchical data storage structures

ISB

20

# Intro to DC from APIs

ISB

21

## DC from APIs: A Simple Weather API...

- https://darksky.net/dev

# Dark Sky API

The easiest, most advanced, weather API on the web.

TRY FOR FREE

Easy to Use          Weather Conditions          Advanced Data

22

**DC from APIs: A Simple Weather API…**

- Who in business might want to use a weather API? Why?

- Examine what the documentation says:

- What the API gives in terms of output data fields

- Which among those fields are required vs optional

- How to construct a query for those fields

- Some sample output

- Let's run the queries we've built and examine output
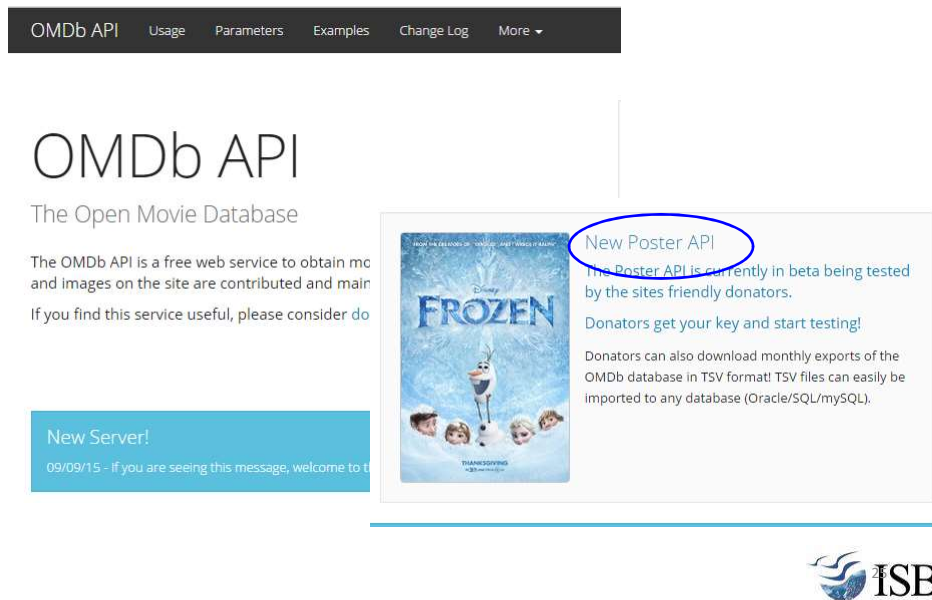
ISB

23

**Recapping our Weather API journey**

- What is an API? How to connect to one?

- How to know what all data fields are available from an API?

- How to know what the API query construction should be like?

- How to know the *pricing* of data and/or services?

- What *after* getting the data from an API?

ISB

24

## Data from APIs: A Simple Movies DB API…

- http://www.omdbapi.com/



25

## Data from APIs: A Simple Movies DB API…

- Goto http://www.omdbapi.com/ and explore the page a bit …

- Qs to ponder upon:

- How is the API key delivered? What are the usage restrictions?

- How does this site make money?

- What does their legal section say?

- Open the file '' and browse through the results of an API query on omdb

26

# A Mini Use-case for DC from APIs



27

---

## Setting the Context: *Dynamic* Digital Advertising

- Consumers dislike advertisements. *Reach* alone isn't enough anymore ..
- Marketers have tried many tips and tricks to somehow, anyhow elicit some response, some engagement…

- And one truism that's emerging is that *Relevance* trumps Reach.

- How might media firms / *publishers* sell "Relevance"?
  - Relevance implies alignment with customer's needs/preferences/ interests etc.
  - Hope is a Customer's demographic and web-surfing profile correlates with preferences/interests etc.
  - Challenge is to (micro-)segment population → for Targeting, Re-targeting
  - If the segment is fine enough, personalized ads created → served to right people at the right time

28

## A Simple *Dynamic* Advertising Campaign

- At its simplest, dynamic ad campaigns require **data signals** (from APIs), a set of *event rules* or *"trigger" conditions* and ad copy/messaging to go along.
- Consider this illustrative example…

| Brands | Some data signals received from APIs… | | | | | |
|---|---|---|---|---|---|---|
| | Weather | Profile | Sports | Entertainment | Geography | Social |
| Audi | ✓ | ✓ | ✓ | ✓ | | ✓ |
| Lipitor | | ✓ | ✓ | ✓ | ✓ | |
| L'Oreal | ✓ | ✓ | | ✓ | ✓ | ✓ |
| Wells fargo | | ✓ | | | ✓ | ✓ |
| Starbucks | ✓ | | | | | ✓ |

Discretion and judgment are important. Some sports (say, football) gel better with some brands (Bud) than others.

- *What* events might trigger *which* ad copies to get served to *whom*?

ISB

29

## A Simple Dynamic Advertising Campaign

- Again, illustrative only…

| Data Signal(s) | Trigger Condition | Message | Images |
|---|---|---|---|
| Weather, City | Weather = Rainy, City = Seattle | Hey, Seattle, it's been pouring. Time to hit the Gym. | Animated Gym scene, treadmill, Nike jogging shoe |
| Weather, City | Weather = Cold, City = Seattle | Hey, Seattle, it's freezing. Warm up with a workout. | Animated image of huddled person vs person on a treadmill, Nike jogging shoe. |
| Weather, City | Weather = Hot, City = Seattle | Hey, Seattle, it's a glorious sunny day. | Golf range, Golfer in Nike Golf shoes. |

- Marketers write "event rules" and creatives build messaging copy (including the default copy) in tandem.

- Modularization of ad-copy into elements that can be mixed and matched to dynamically create ads is a reality today.

ISB

30

## Digital Advertising: Some Data Signals and Events

- Some common data sources and event "triggers":
- 1. **Profile data** (aka basic demographics)

  E.g., in SUV ads, show Merc only in some geocodes; show ruggedness to males but emphasize space or safety to females, etc.

- 2. **CRM data** - Purchase history, site browsing history, brand loyalty, payment methods etc.
  E.g., showing gluten-free food ads only to gluten allergic people.

- 3. **Environmental** data - real-time weather, temp, geolocation, date & time...

  E.g., "TGIF, go home, have a Kingfisher, welcome the weekend..."; "Chilly days call for [Starbucks] Cappucino…"

- 4. **Social media data** - trends, likes and preferences, topics discussed of late

  E.g., "Dravid endorses AkshayaPatra" to Cricket enthusiasts/ people who've liked/shared some recent Dravid related news.

ISB

31

## Digital Advertising: Data Signals and Events

- 5. **Real-time events** - usually big sports or entertainment events around which some clever messaging, copy or campaigning can be built.
- E.g., "Go, Mumbai Indians!", "Buy your Valentine a Swatch."

- 6. **Site/cookie data** - first party cookies on consumer websites
- E.g., "Traveling? Try the new VIP series 6 suitcases…"

- 7. **Search data** - SEO etc. But advertising based on this is harder now that Google has moved to secure search and hides the referrer URL's search term...

- 8. **Contextual data** - media section based. E.g., is user browsing finance section? sports? lifestyle? health?

ISB

32

Session Wrap-up

## Session Wrap-up

- What did we cover this session – some salient points?

- What other APIs are out there – what would you like to have seen?

- Ready for reasonably simple homework assignments on this topic?

- Any other Qs or comment?