



Clase 3

Webscraping y API's

Aproximación a las políticas públicas desde los datos | UC | 09 de junio, 2023

👤 José D. Conejeros | ✉️ jdconejeros@uc.cl

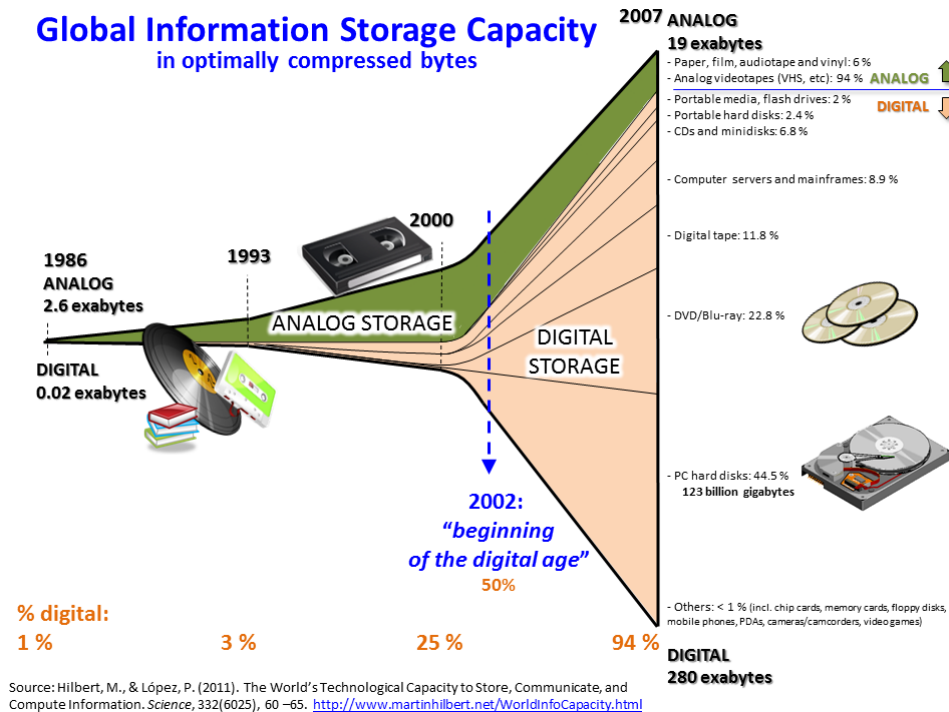
Guía

1. **Extracción de datos vía webscraping**
2. Extracción de datos vía uso de APIS's

La información en internet

La información en Internet es abundante y está al alcance de la mano, pero hay desafíos en:

- Cómo recopilar la información → Webscraping/Api's
- Cómo analizarla la indormación → `tidydata` (lo veran a lo largo del diplomado)



¿Qué es el webscraping?

Es el proceso de extracción de datos estructurados, semi estructurados y no estructurados (HTML) de forma automática desde un sitio web.

- Datos estructurados: datos que tienen un modelo definido o provienen de un campo determinado en un registro. Por ejemplo: el clásico marco de datos con filas y columnas.
- Datos semi estructurados: datos que no tienen formatos fijos, pero contienen atributos o etiquetas. Por ejemplo: correos electrónicos, **páginas HTML**
- Datos no estructurados: datos que no tienen un modelo predefinido o no están organizados de alguna manera. Por ejemplo: videos, fotografías, audio, texto.

Buscamos transformar los datos no estructurados a una estructura conocida del tipo bidimensional: filas y columnas, **aunque no siempre esto es lo más conveniente**. ¿De qué depende?

Algunos ejemplos de uso de webscraping

- Clasificar productos o servicios para crear motores de recomendación, también nos sirve para obtener datos de texto, como en Wikipedia o para hacer sistemas basados en el Procesamiento del Lenguaje Natural.
- Generar datos a partir de etiquetas de imágenes, de sitios web como Google, Flickr, etc. para entrenar modelos de clasificación de imágenes.
- Consolidar los datos de las redes sociales: Facebook y Twitter, para realizar análisis de sentimientos, opiniones o *topic modeling*.
- Extraer comentarios de usuarios y sitios de comercio electrónico como Alibaba, Amazon o Walmart, etc.
- ¡Todas las que puedan imaginar!

Formas de extraer información: DOM

- **Análisis DOM (Documento Objeto Modelo):** interfaz de plataforma que proporciona un conjunto estándar de objetos para representar documentos en HTML y XML.
- El DOM permite el acceso dinámico a través de la programación para ingresar, añadir y cambiar dinámicamente el contenido estructurado en documentos a través de lenguajes como *JavaScript*.
- Los objetos DOM modelan tanto la ventana del navegador como el historial, el documento o la página web, y todos los elementos que la propia página puede tener: párrafos, divisiones, tablas, formularios y sus campos, etc.
- A través del DOM se puede acceder, mediante *Javascript*, a cualquiera de estos elementos, es decir, a sus correspondientes objetos para alterar sus propiedades o invocar sus métodos.

Formas de extraer información: DOM

The image shows a web browser displaying a dashboard titled "Detector de riesgos y desastres". The dashboard includes a header, a main content area with a title "Alertas de riesgos y desastres en prensa local", and a data visualization section with four donut charts. The charts show: "Noticias revisadas" (6.885), "Alertas detectadas /30 días" (170), "Regiones con alertas" (15), and a bottom row with "refieren a Ecología" (4), "mencionan Peligros" (65), "mencionan Desastres" (26), and "mencionan Eventos" (80). A sidebar at the bottom allows selecting regions and areas. A Chrome DevTools DOM inspector is overlaid on the right, showing the HTML structure of the page, including an embedded content area.

Alertas de riesgos y desastres en prensa local

Métrica	Valor
Noticias revisadas	6.885
Alertas detectadas /30 días	170
Regiones con alertas	15
refieren a Ecología	4
mencionan Peligros	65
mencionan Desastres	26
mencionan Eventos	80

DOM Inspector (HTML Structure):

```
<section id="comp-17c44d87" class="comp-17c44d87 CohWsy wixui-column-strip">...</section>
<div id="comp-17c44d9d1" class="KcpHe0 tz5f0K comp-17c44d9d1 wixui-rich-text" data-testid="richTextElement" data-angle="0" data-angle-style-location="style" style="visibility: inherit;" data-screen-in-hide="done">...</div>
<div id="comp-17c44d91" class="hFQZvn comp-17c44d91">
  <div class="nTiihl wixui-box" data-testid="container-bg">...</div>
  <div data-mesh-id="comp-17c44d9iinlineContent" data-testid="inline-content" class="...>...</div>
</div>
<div id="comp-17c44d9s" class="hFQZvn comp-17c44d9s">...</div>
<div id="comp-17c44d9m" class="comp-17c44d9m _xg6_p">
  <wix-iframe data-src="...">
    <div class="SMuTia">
      <iframe class="wuksD5" title="Embedded Content" name="htmlComp-iframe" width="100%" height="100%" data-src src="https://www.monitorsocial-cl.filesusr.com/html/657f14_4cfc914_.html" data-gtm-yt-inspected-4="true">...</iframe>
    </div>
  </wix-iframe>
</div>
<div id="comp-17c44d95fn" class="hFQZvn comp-17c44d95fn">...</div>
<div id="comp-17c44d9jh2" class="hFQZvn comp-17c44d9jh2">...</div>
<div id="comp-17c44d9g9" class="hFQZvn comp-17c44d9g9">...</div>
<div id="comp-17c44d9v" class="IA2mX comp-17c44d9v">...</div>
<div id="comp-17c44d9x3" class="KcpHe0 tz5f0K comp-17c44d9x3 wixui-rich-text" data-testid="richTextElement" data-angle="0" data-angle-style-location="style" style="visibility: inherit;" data-screen-in-hide="done">...</div>
<div id="comp-17c44d9z1" class="KcpHe0 tz5f0K comp-17c44d9z1 wixui-rich-text" data-testid="richTextElement" data-angle="0" data-angle-style-location="style" style="visibility: inherit;" data-screen-in-hide="done">...</div>
<div id="comp-17c44da02" class="KcpHe0 tz5f0K comp-17c44da02 wixui-rich-text" data-testid="richTextElement" data-angle="0" data-angle-style-location="style" style="visibility: inherit;" data-screen-in-hide="done">...</div>
<div id="comp-17c44da21" class="KcpHe0 tz5f0K comp-17c44da21 wixui-rich-text" data-testid="richTextElement" data-angle="0" data-angle-style-location="style" style="visibility: inherit;" data-screen-in-hide="done">...</div>
```

Monitor Social

Formas de extraer información: DOM

Podemos extraer datos del código fuente del sitio web, con un analizador html y expresiones regulares:

```
# Paquete de R que permite la extracción  
install.packages("rvest")  
library(rvest)
```

Puede ver la documentación de la librería dado [click aquí](#)

- Ojo con las condiciones de servicio: generalmente en la parte inferior del sitio web
- Verificar que los datos que se extraen sean de acceso público
- Cuidado con recopilar información personal
- Revisar los derechos de autor

Ojo a la hora de hacer extracciones

1. Respete los deseos del sitio de alojamiento:

Compruebe si existe una API o si los datos están disponibles para descargar

Tenga en cuenta de dónde provienen los datos y dé el crédito respectivo: respete los derechos de autor si desea publicar los datos

Algunos sitios web no permiten extracciones de datos en el archivo `robots.txt` (lo veremos)

2. Limite sus extracciones:

Espere uno o dos segundos después de cada extracción: los servidores web identifican las extracciones de datos

Extraiga solo lo que necesita, y solo una vez (por ejemplo, almacene el archivo html y luego analícelo)

Grandes extracciones puede demandar mucho al servidor. Esto puede acarrear problemas y finalmente las personas a cargo de la plataforma decidan no disponer de los datos de interés.

robots.txt

Es un archivo presente en la mayoría de los sitios web. Contiene el llamado estándar de exclusión de robot que son una serie de instrucciones especialmente dirigidas a los programas que buscan indexar el contenido de estas páginas (por ejemplo el bot de Google que “guarda” las nuevas páginas que se crean en Internet).

- Este método se utiliza para evitar que ciertos bots que analizan sitios de Internet añadan información “innecesaria” a los resultados de las búsquedas.
- Un archivo [robots.txt](#) en una página web funcionará como una petición para que ciertos bots ignoren archivos o directorios específicos en su búsqueda.
- Se aconseja siempre revisar este archivo antes de iniciar el proceso de extracción.

```
User-agent: *
Disallow: /search
Allow: /search/about
Allow: /search/static
Allow: /search/howsearchworks
Disallow: /sdch
Disallow: /groups
Disallow: /index.html?
Disallow: /?
Allow: /?hl=
```

[robots.txt](#)

robots.txt

```
User-agent: *  
Allow: /  
Disallow: /*/alt-*.html  
Disallow: /*/aux-*.html  
Crawl-Delay: 2
```

robots.txt

User-agent: * indica a quién se aplican estas instrucciones. A veces aparece el nombre de algún robot en particular. En este caso, el asterisco explicita que estas instrucciones aplican a todos.

Allow: / indica cuáles son las rutas dentro del sitio desde las que es posible extraer datos. En este caso la barra "/" indica que es posible hacerlo de todas.

Disallow: /*/alt-*.html y **Disallow: /*/aux-*.html** explicitan aquellas rutas dentro del sitio web desde las que no se da autorización para hacer una extracción de forma automatizada. En este caso, son todas aquellas páginas cuya url incluye alt- o aux-. Es decir, acá se está poniendo una restricción a la autorización general que nos dieron en la línea anterior.

Crawl-Delay: 2 señala la cantidad de segundos que es necesario esperar entre cada petición (request) de datos al sitio. Esto es importante tenerlo en cuenta cuando extraemos datos de varias páginas dentro de un sitio web al mismo tiempo, ya que esos segundos de espera permiten que el funcionamiento del sitio no se vea afectado.

Flujo de trabajo en un Webscraping

- a. Investigar y comprender la estructura del sitio web.
- b. Elegir una estrategia de extracción: objetivos y outcomes de interés.
- c. Cree una maqueta del código: extraiga, evalúe y ajuste.
- d. Generalice su maqueta para automatizar la extracción: uso de funciones, bucles for y depuración.
- e. Aplique una exploración y limpieza de su extracción: debe considerar que esto toma bastante tiempo.

Antes de hacer una extracción

Página web=archivo de texto etiquetado

Lenguaje de marcado de hipertexto (HTML): estándar oculto detrás de cada sitio web HTML es texto con estructura marcada, definida por etiquetas:

```
<html>
<head>
  <title>Mi website</title>
</head>
<body>
  <h1 id='first'>Título de mi web</h1>
  <p>Texto & <b>Texto en negrita.</b>
  <img src='myimg.png' width='100' height
  <h2 class="encabezado-lista">Título lis
  <ol>
    <li>Un primer elemento</li>
    <li>Un segundo elemento</li>
    <li>Un tercer elemento</li>
  </ol>
  <h2 class="encabezado-parrafo">Otro títul
</body>
```

Título de mi web

Texto & **Texto en negrita.**



Título lista

1. Un primer elemento
2. Un segundo elemento
3. Un tercer elemento

Otro título

Antes de hacer una extracción

Lo que hacemos al realizar webscraping es importar el código html de un sitio web a nuestro computador (usando R, por ejemplo) y extraer aquellas partes específicas que nos interesan.

- El contenido de un archivo html se encuentra etiquetado, podemos utilizar esas mismas etiquetas para especificar cuáles son las partes de la página que buscamos extraer.
- Las etiquetas sirven para distinguir elementos, clases y, en algunos casos, el id de sus distintos componentes.
- Las etiquetas pueden tener atributos con nombre que se parecen a `name1='value1'` `name2='value2'`. Dos de los atributos más importantes son `id` y `class`, que se utilizan junto con el archivo CSS (hojas de estilo en cascada) lo que permite controlar la apariencia visual de la página.
- Entender la diferencia entre un “elemento” del documento html y una “clase” nos va ayuda a precisar mejor cuál es el contenido exacto que queremos de una página.

Antes de hacer una extracción

Hay más de 100 elementos HTML. Algunas etiquetas comunes:

Elementos del documento: `<head>`, `<body>`, `<footer>`, ...

Componentes del documento: `<title>`, `<h1>`, `<div>`, ...

Estilo de texto: ``, `<i>`, ``, ...

hipervínculos: `<a>`

El web scraping es posible porque la mayoría de las páginas que contienen datos tienen una estructura consistente.

Puede revisar más sobre etiquetado html en: [MDN web docs](#)

Una herramienta útil para identificar selectores: [selectorGadget](#)

Webscraping de páginas dinámicas

Las extracciones funcionan bastante bien con páginas web estáticas donde el sitio web y sus componentes no interactúan con el usuario. Sin embargo, hay páginas web que son dinámicas pues:

- Cambian a medida que el usuario interactúa con él
- Tienen elementos reactivos (por ejemplo, menús desplegables)

En este escenario será necesario controlar el navegador para realizar extracciones en este tipo de páginas web. Algunos paquetes útiles para esto: [RSelenium](#) o [chromote](#).

Esto excede los objetivos del curso, pero pueden encontrar guías de uso:

[Web RSelenium](#)

[Guía general](#)

[Video de instalación](#)

[Video guía de uso](#)

En simple, se iniciará una sesión del navegador a través de R y toda la comunicación se enrutará a partir de dicha sesión.

Guía

1. Extracción de datos vía webscraping
2. **Extracción de datos vía uso de APIS's**

Application Program Interfaces (APIs)

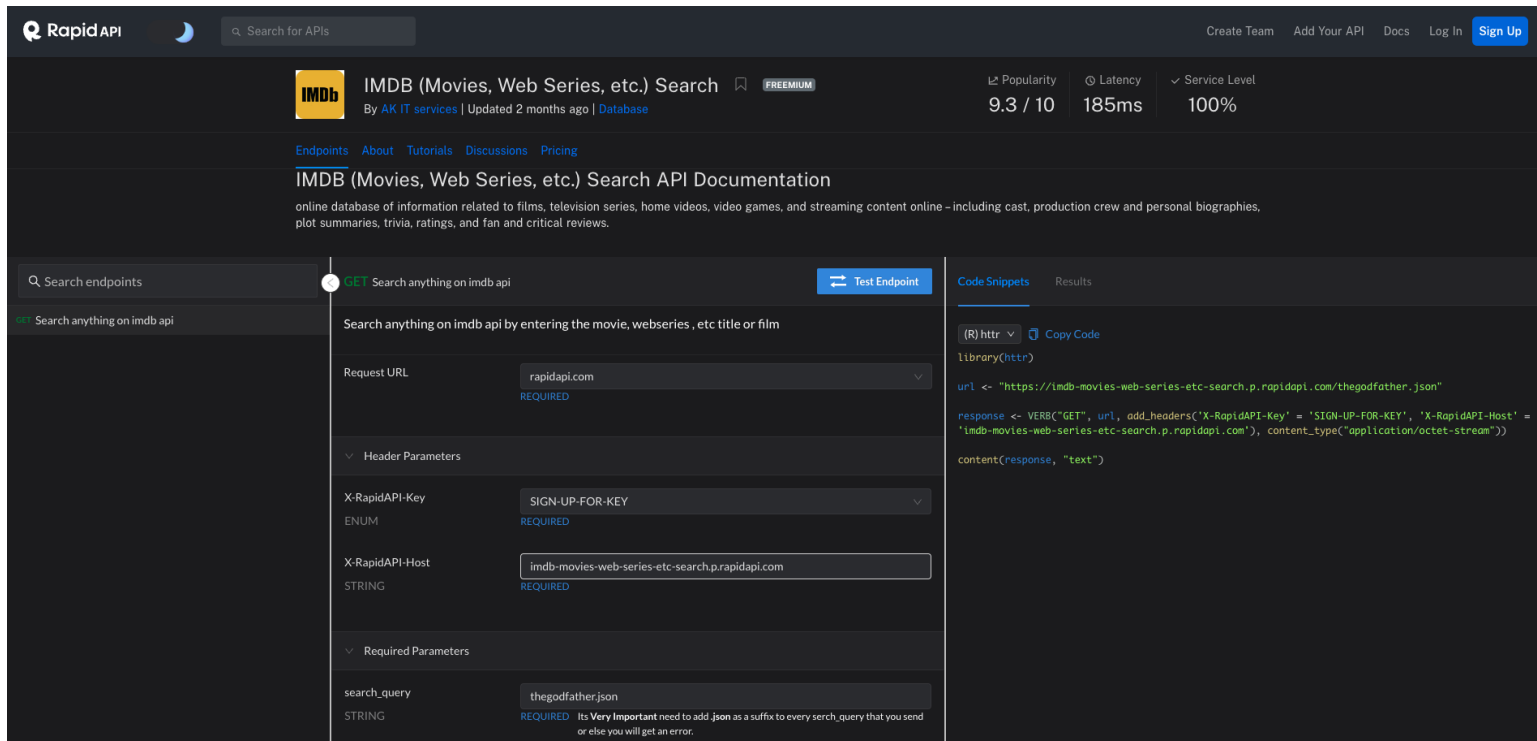
APIs (interfaz de programación de aplicaciones): es un conjunto de definiciones y protocolos utilizados para desarrollar e integrar programas informáticos de aplicación.

- Es un código que indica a las aplicaciones cómo pueden comunicarse entre sí. Las API permiten que sus productos y servicios se comuniquen con otros sin necesidad de saber cómo se implementan.
- Sitios web como Facebook, Twitter, LinkedIn, entre otros, ofrecen una API pública y/o privada a la que se puede acceder mediante programación para recuperar datos en el formato deseado. No todas las páginas webs tienen API's disponibles.
- Para que las APIs públicas sean utilizadas muchos sitios solo permiten a usuarios autorizados (por ejemplo aquellos que tienen una cuenta en la plataforma).
- Siempre que sea posible debe utilizar la API, ya que, por lo general, le brindará datos más confiables.

CUIDADO: revisar los términos y condiciones de extracción, pues existen cargos para tarjetas de crédito si se superan los límites de extracción pública.

Formas de extraer información: APIs

- [httr](#): herramienta general para hacer extracciones de información desde APIs en internet.
- [rapidapi](#): web con varias APIs disponibles para la extracción de datos. Utilizan varios códigos de extracción.



The screenshot displays the RapidAPI interface for the IMDB (Movies, Web Series, etc.) Search API. The page header includes the RapidAPI logo, a search bar, and navigation links like 'Create Team', 'Add Your API', 'Docs', 'Log In', and 'Sign Up'. The API details section shows the IMDB logo, the API name 'IMDB (Movies, Web Series, etc.) Search', and its popularity (9.3 / 10), latency (185ms), and service level (100%). Below this, there are links for 'Endpoints', 'About', 'Tutorials', 'Discussions', and 'Pricing'. The main content area is titled 'IMDB (Movies, Web Series, etc.) Search API Documentation' and describes it as an 'online database of information related to films, television series, home videos, video games, and streaming content online - including cast, production crew and personal biographies, plot summaries, trivia, ratings, and fan and critical reviews.' The interface is split into three main sections: a search bar on the left, a configuration panel in the middle, and a code editor on the right. The search bar contains 'Search anything on imdb api'. The configuration panel shows the 'Request URL' as 'rapidapi.com', 'Header Parameters' including 'X-RapidAPI-Key' (set to 'SIGN-UP-FOR-KEY'), 'X-RapidAPI-Host' (set to 'imdb-movies-web-series-etc-search.p.rapidapi.com'), and 'Required Parameters' including 'search_query' (set to 'thegodfather.json'). The code editor shows an R snippet using the httr library to make a GET request to the API endpoint and print the response content.

API IMDB (Movies, Web Series, etc.) Search

Librerías R

Existen librerías en R con los cuales se pueden trabajar extracciones a partir de API's. Algunos ejemplos:

- [spotifyr](#): para extraer funciones de audio de pistas y otra información de la API web de Spotify de forma masiva
- [WikipediR](#): para extraer información de wikipedia
- [rtimes](#): permite extraer noticias del New York Times
- [ubeR](#): extracción de datos desde UBER
- [Rlinkedin](#): extracción de datos desde LinkedIn
- [googleAnalyticsR](#): extrae información desde Google Analytics
- [Rbnb](#): acceso a la API de Airbnb
- [ggmap](#): acceso a mapas de google
- [academictwitterR](#): recopilar tweets desde el punto final de la API v2 para el seguimiento del producto de investigación académica
- [rtweet](#): para extraer información desde la API de twitter

Un largo listado se encuentra aquí: [Lista de paquetes que conectan con API's](#).

Notas: Antes de realizar una extracción de datos es siempre bueno verificar si hay un paquete o API disponible.

Aspectos relevantes a la hora de hacer extracciones

- Crearse una cuenta para acceder a la API.
- Al usar API, lea la documentación
 - ¿Hay una opción de descarga por lotes?
 - ¿Hay algún límite de extracción?
 - ¿Desde qué punto se deben pagar los datos?
 - ¿Puedes compartir los datos?

Referencias útiles

Webscraping:

- [Minería de datos web](#)
- [Web scraping R for data science](#)
- [Data collection: Web \(screen\) scraping](#)

Riva Quiroga, "[Introducción al web scraping usando R](#)", Programming Historian en español 6 (2022), <https://doi.org/10.46430/phes0061>.

API's:

- [Using Web APIs from R](#)
- [Accessing Web APIs](#)

Curso de Brian Cooksey [An introduction to APIs](#)



Clase 3

Webscraping y API's

09 de junio, 2023

 **José D. Conejeros** |  jdconejeros@uc.cl |  JDConejeros