

# Methods for unstructured data

## Lecture 2: Web scraping

---

Helge Liebert

# Roadmap and goals

- Use processing of web data to introduce various ideas along the way.
- Accessing data on the web: APIs.
- Gathering (semi-structured) web data and transforming it into structured data (*“web scraping”*).
- Give you a general understanding how web scrapers work.
- Foster an understanding of the development process.
- Point you towards the necessary tools so you can write your own.

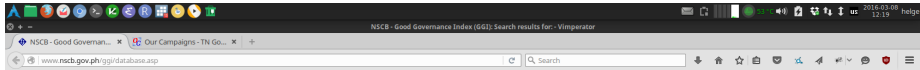
1. Introduction
2. Static and dynamic websites
3. Application Programming Interfaces
4. The Document Object Model
5. Examples
6. Assignment

# Guiding problem

- How to turn unstructured into structured data?
- Consider a situation where
  - You want to get data from the internet.
  - The data is in unstructured/semi-structured form.
  - Possibly embedded in a website.
  - You want to transform it into a differently structured format for further use.
  - You need to filter the available information.



...or this



#### Good Governance Index

##### Database

##### Municipality

Search for municipality

Go

or choose from the list

- NCR - Metro Manila, Navotas
- Pr NCR - Metro Manila, Pateros
- NCR - Metro Manila, San Juan
- CAR - Abra, Bangued
- CAR - Abra, Boliney
- CAR - Abra, Burrey
- CAR - Abra, Bucloc
- SE CAR - Abra, Dagupan
- CAR - Abra, Danguas
- fo CAR - Abra, Dolores
- SE CAR - Abra, Lala
- CAR - Abra, Lacub
- SE CAR - Abra, Lagangilang
- CAR - Abra, Lagayan
- CAR - Abra, Langiden
- CAR - Abra, Uman-Baay
- CAR - Abra, Luba
- CAR - Abra, Malibcong
- CAR - Abra, Manabo
- CAR - Abra, Pinarubia

##### Good Governance Index

Main Page

Database

Technical Notes

Press Release

E-mail the webmaster

Terms of Use

Home - Top of Page

PHILIPPINE STATISTICS AUTHORITY  
PSA CVEA Bldg., East Avenue Dillman, Quezon City, Philippines  
Telephone no. 462-6600 loc. 839; E-mail: info@psa.gov.ph

INSERT

http://www.nscb.gov.ph/ggi/database.asp [1/2] All

Enging - Chromium

500 - Internal server error...

Slides.tex (-)Dropbox/Pres...

Inbox - hyperboloidcount...


Funktion mit Style - Neue...

NSCB - Good Governance L...

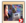
# ...and this

NSCB - Good Governance x Our Campaigns - Candidates x

www.nscb.gov.ph/ggi/details\_prov.asp



REPUBLIC OF THE PHILIPPINES  
**PHILIPPINE STATISTICS AUTHORITY**  
SOLID • RESPONSIVE • WORLD-CLASS


**Good Governance Index**

**GGI Indicators for the province of Misamis Oriental**

Indicator	2005		2008	
	Value	Rank in	Value	Rank in
<b>Economic Governance Index</b>	<b>107.08</b>	<b>37</b>	<b>86.12</b>	<b>71</b>
Total Financial Resources Generated per Capita (in Millions)	171.05	45	233.43	50
Per Capita Tax and Non-Tax (in Millions)	71.66	1	73.87	7
Per Capita Total Deposits	122.61	4	N/A	N/A
Per Capita Expenditure on Social Services	212.80	22	90.22	67
Unemployment Rate	125.56	47	125.56	47
Underemployment Rate	55.31	53	55.31	53
Inflation Rate	66.37	14	37.46	37
Poverty Gap	87.56	34	72.64	42
Poverty Incidence	84.90	27	88.10	31
<b>Political Governance Index</b>	<b>103.24</b>	<b>13</b>	<b>99.86</b>	<b>18</b>
Crime Solution Efficiency Rate	103.91	8	96.76	23
Voter's Turnout Rate	102.56	31	102.56	31
<b>Administrative Governance Index</b>	<b>128.00</b>	<b>23</b>	<b>154.38</b>	<b>13</b>
<b>A. Education Index</b>	<b>80.79</b>	<b>68</b>	<b>107.25</b>	<b>52</b>
Elementary Teacher to Pupil Ratio	100.07	55	96.93	54
High School Teacher to Pupil Ratio	100.93	45	99.86	53
No. of Public Elem. Schools per 1000 School Age Population	56.02	85	96.64	45
No. of Public High Schools per 1000 School Age Population	51.16	70	118.54	60
Enrollment in Gov't. Elem. Sch. per 1000 School Age Population	57.08	75	95.78	9
Enrollment in Gov't. HS. Sch. per 1000 School Age Population	61.13	63	140.90	8
Cohort Survival Rate (Elementary)	93.57	46	83.99	49
Cohort Survival Rate (High School)	79.85	32	85.29	48
Elementary Pupil - Classroom Ratio	102.75	58	103.88	54
High School Pupil - Classroom Ratio	105.37	61	140.74	53
<b>B. Health Index</b>	<b>229.66</b>	<b>19</b>	<b>212.88</b>	<b>21</b>
Total Health Personnel per 1000 Population	193.11	19	170.55	23
% Birth less the 2500 grams	478.47	12	393.89	21
% of Household with access to safe water	109.09	48	131.13	1
% Barangay Health Station per 100,000 Population	137.99	30	151.93	21
<b>C. Power and ICT Index</b>	<b>76.26</b>	<b>35</b>	<b>144.00</b>	<b>1</b>
Power Index	124.83	1	124.83	1
ICT Index (Telephone Density per 1000 Population)	27.89	41	157.44	1
<b>Good Governance Index</b>	<b>116.27</b>	<b>33</b>	<b>117.31</b>	<b>42</b>

1 = highest/best, out of 79 provinces

**Good Governance Index**  
[Main Page](#)  
[Database](#)  
[Technical Notes](#)  
[Press Release](#)

...to this.

master.csv - LibreOffice Calc																							
File Edit View Insert Format Tools Data Window Help																							
Liberation Sans 10																							
A1																							
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R						
id	name	jahr	land	provinz	gruppe	popcat	pop	ehe	lebgeb	lebgebo	lebgebot	todgebo	todgebot	geb	gebo	todgeb	todgebshare						
1	Aachen	1927	Preussen	Rheinprovinz	A	pop >= 100,000	156360	1423	2753	2366	387	80	7	2840	2446	87	3.0633902						
2	Aachen	1928	Preussen	Rheinprovinz	A	pop >= 100,000	150991	1446	2723	2260	463	66	17	2806	2326	83	2.9579473						
3	Aachen	1929	Preussen	Rheinprovinz	A	pop >= 100,000	155542	1558	2635	2192	443	49	23	2707	2241	72	2.659771						
4	Aachen	1930	Preussen	Rheinprovinz	A	pop >= 100,000	154634	1501	2635	2189	446	46	22	2703	2235	68	2.5157232						
5	Aachen	1931	Preussen	Rheinprovinz	A	pop >= 100,000	154400	1408	2429	2030	399	53	22	2504	2083	75	2.9952078						
6	Aachen	1932	Preussen	Rheinprovinz	A	pop >= 100,000	153834	1426	2305	1968	337	46	20	2371	2014	66	2.7836356						
7	Aachen	1933	Preussen	Rheinprovinz	A	pop >= 100,000	162990	1616	2371	2028	343	54	23	2448	2082	77	3.1454248						
8	Aachen	1934	Preussen	Rheinprovinz	A	pop >= 100,000	163839	1942	2841	2451	480	76	8	3025	2527	84	2.7768595						
9	Aachen	1935	Preussen	Rheinprovinz	A	pop >= 100,000	164180	1570	3048	2516	532	54	21	3123	2570	75	2.4015369						
10	Aachen	1936	Preussen	Rheinprovinz	A	pop >= 100,000	164105	1502	3012	2382	630	43	25	3080	2425	68	2.2077923						
12	Ahlen	1927	Preussen	Westfalen	D	15,000 <= pop < 30,000	23956	239	606	596	10	27	1	634	623	28	4.4164038						
13	Ahlen	1928	Preussen	Westfalen	D	15,000 <= pop < 30,000	24703	266	625	621	4	18	0	643	639	18	2.7993779						
14	Ahlen	1929	Preussen	Westfalen	D	15,000 <= pop < 30,000	25043	227	624	609	15	26	0	650	635	26	4						
15	Ahlen	1930	Preussen	Westfalen	D	15,000 <= pop < 30,000	25226	202	568	548	20	25	1	594	573	26	4.3771043						
16	Ahlen	1931	Preussen	Westfalen	D	15,000 <= pop < 30,000	25373	191	508	483	25	34	0	542	517	34	6.2730627						
17	Ahlen	1932	Preussen	Westfalen	D	15,000 <= pop < 30,000	25549	200	550	514	36	13	1	564	527	14	2.4822695						
18	Ahlen	1933	Preussen	Westfalen	D	15,000 <= pop < 30,000	25153	301	478	437	41	15	2	495	452	17	3.4343433						
19	Ahlen	1934	Preussen	Westfalen	D	15,000 <= pop < 30,000	25700	283	658	606	52	19	2	679	625	21	3.0927835						
20	Ahlen	1935	Preussen	Westfalen	D	15,000 <= pop < 30,000	25937	212	661	621	40	17	0	678	638	17	2.5073745						
21	Ahlen	1936	Preussen	Westfalen	D	15,000 <= pop < 30,000	26104	228	673	625	48	11	1	685	636	12	1.7518249						
22	Allenstein	1927	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	39315	241	892	859	33	14	3	909	873	17	1.870187						
23	Allenstein	1928	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	39232	293	943	898	45	24	2	969	922	26	2.6831784						
24	Allenstein	1929	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	39114	267	959	900	56	21	3	924	855	25	2.5406504						
25	Allenstein	1930	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	39345	283	862	810	52	14	5	881	824	19	2.1566403						
26	Allenstein	1931	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	39876	301	884	836	48	17	4	905	853	21	2.320442						
27	Allenstein	1932	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	40078	323	837	796	41	15	1	853	811	16	1.8757327						
28	Allenstein	1933	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	43079	340	890	854	36	14	4	908	868	18	1.9823788						
29	Allenstein	1934	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	43506	405	1062	1015	47	28	4	1094	1043	32	2.9250457						
30	Allenstein	1935	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	43881	363	1118	1051	67	21	4	1138	1067	20	1.7574693						
31	Allenstein	1936	Preussen	Ostpreussen	C	30,000 <= pop < 50,000	44566	364	1150	1037	82	16	5	1145	1058	26	2.2707424						
32	Altena	1927	Preussen	Westfalen	D	15,000 <= pop < 30,000	15931	131	240	226	14	5	0	245	231	5	2.0408163						
33	Altena	1928	Preussen	Westfalen	D	15,000 <= pop < 30,000	16333	163	244	229	15	11	2	257	240	13	5.0583658						
34	Altena	1929	Preussen	Westfalen	D	15,000 <= pop < 30,000	16464	139	257	248	9	6	1	264	254	7	2.6515152						
35	Altena	1930	Preussen	Westfalen	D	15,000 <= pop < 30,000	16498	139	242	220	22	16	1	259	236	17	5.6537064						
36	Altena	1931	Preussen	Westfalen	D	15,000 <= pop < 30,000	16407	138	196	186	10	5	2	203	191	7	3.4482758						
37	Altena	1932	Preussen	Westfalen	D	15,000 <= pop < 30,000	16115	122	192	176	16	9	0	201	185	9	4.477612						
38	Altena	1933	Preussen	Westfalen	D	15,000 <= pop < 30,000	16133	138	189	174	15	6	1	196	180	7	3.5712695						
39	Altena	1934	Preussen	Westfalen	D	15,000 <= pop < 30,000	16246	162	258	223	35	8	2	268	231	10	3.7313433						

Find

Find All Match Case

Sheet 1 of 1

Default

Sum=0

110%



# Goal: Automation

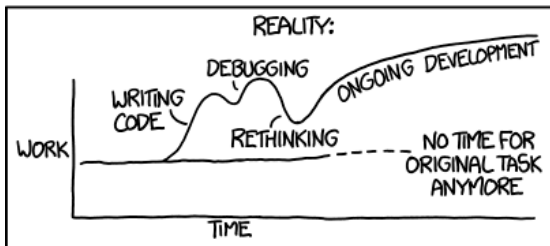
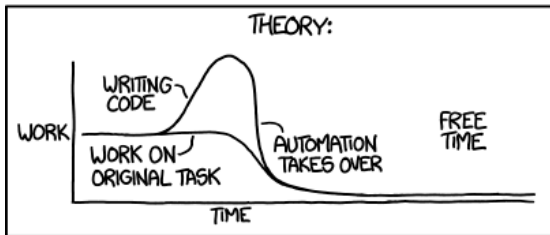
- Digitalization offers exciting data for research. But: *Data is messy*.
- Gathering or processing data often involves repetitive manual tasks.
- Disadvantages:
  - Manual tasks are often not well documented or reproducible ex post.
  - Manual work is frustrating and a huge time-sink.
  - Manual work may not be feasible with large data.

## ➡ Automation helps!

- Frees you to engage in other work.
- You learn new things.
- Should you encounter the same class of problem in the future, you already have a solution at hand.

# Automation

"I SPEND A LOT OF TIME ON THIS TASK.  
I SHOULD WRITE A PROGRAM AUTOMATING IT!"



## Getting started—things to consider before you begin

- Pick up the phone and try to get the data directly.
- Search if somebody has already faced the same or a similar problem.
- Does the site or service provide an API that you can access directly?
- Is there a wrapper for it?
- Is the website only online for a limited time? Do you want an original snapshot as a backup? Is it more convenient to filter your data offline?

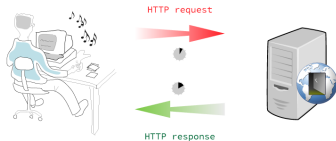
## Static and dynamic websites

---

# Static vs. dynamic websites

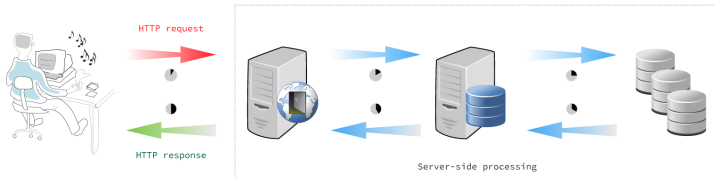
Scheme A

## Static Website



Scheme B

## Dynamic Website



## Save an offline copy

- Use the shell utilities `wget` or `curl` to download the complete site.
- Also useful if you just want a set of files (e.g. pdf documents) from the same site directory.
- Convenient for static sites of limited size.
- Infeasible for large sites or sites that create content dynamically.

# Examples

- Simple http GET request.  
`wget http://www.google.com`
- Recursively download a website.  
`wget -r http://www.some-site.com/some-subdir/`
- Download all pdfs from a site.  
`wget -r -A.pdf http://url-to-webpage-with-pdfs/`
- Mirror a site offline and convert links for local browsing.  
`wget --mirror -p --convert-links -p ./local-dir  
http://target-website.com`

# Application Programming Interfaces

---



- Data providers often offer Web APIs (*Application Programming Interface*) to access data.
- Allow programmable access to data via a defined set of HTTP messages. Similar to visiting a website: you specify a URL and information is sent to your machine.
- With a website, you receive code interpreted by your browser (HTML, CSS, JavaScript). With an API, you receive data.
- Usually in JSON (*JavaScript Object Notation*) or XML (*Extensible Markup Language*) format.

# Web APIs

- Often just two steps:
  1. Construct the URL query that serves as the API request.
  2. Process the response message the API sends back.
- Examples:
  - <https://api.kivaws.org/v1/loans/newest.html>
  - <https://api.kivaws.org/v1/loans/newest.json>
  - <https://api.kivaws.org/v1/loans/search.json?sector=Agriculture&country=VN>
  - <https://www.theyworkforyou.com/api/getMPs?&key=someapikeyhere&output=js>
- Libraries may offer wrappers for APIs: **WDI**, **wbstats**, **twfy**, **pvsR**, Google Maps, OpenStreetMap/OSRM, ...
- Sometimes it is possible to reverse engineer a site's internal API rather than scraping the HTML.

# The Document Object Model

---

# HTML and the Document Object Model

- Extracting information from the web requires a basic understanding of HTML and the associated Document Object Model (DOM).
- HTML elements provide the structure and content of web pages.
- Consist of `<start>` and `</end>` tags, with content in between.  
`<tagname>Content here</tagname>`
- A page consists of nested elements.
- The `html` element is the outer-most element, nesting the `head` and `body` elements, which in turn have nested elements.
- Nesting structure of elements can be represented by a tree (DOM).

# Document Object Model

- The DOM is a programming interface for HTML and XML documents.
- Provides a structured representation of the document.
- A document as a group of nodes, each node representing a part of the document.
- Allows programmatic access to the tree to change the structure, style and content of the document.
- Connects web pages to scripts or programming languages.

# A simple HTML page

A simple HTML page:

```
<html>
  <head>
    <title>My Web Page</title>
  </head>
  <body>
    <h1>Welcome To My Web Page</h1>
  </body>
</html>
```

How a browser renders this page:

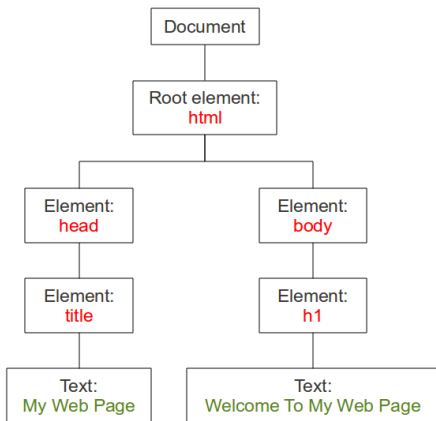


# HTML and the DOM

A simple HTML page:

```
<html>
  <head>
    <title>My Web Page
  </title>
</head>
<body>
  <h1>Welcome To My Web Page
</h1>
</body>
</html>
```

Corresponding node tree:

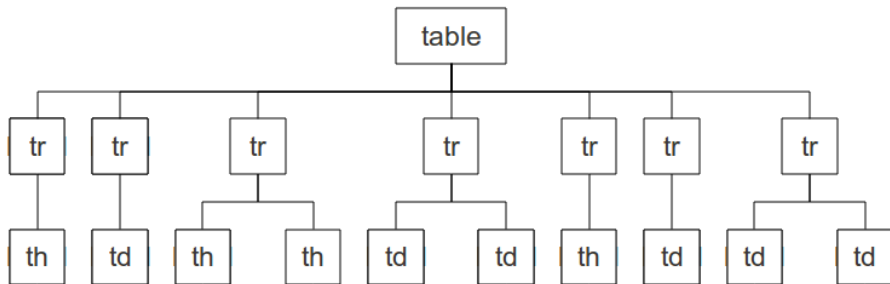


# DOM node trees

- HTML DOM views a document as a tree structure called node tree.
- Everything in an HTML document is a node.
  - The entire document is a document node
  - Every HTML element is an element node
  - Every HTML attribute is an attribute node
  - Text content in the HTML elements is a text node
- Nodes can be accessed through the tree.
- Nodes may be assigned unique id attributes.



## Example: An HTML table element



- Tables are represented by a top-level **table** element.
- The **table** element nests **tr** (*table row*) elements.
- These nest **th** (*table header*) and **td** (*table data*) element cells.

# HTML and the DOM

- HTML tags can have attributes and text content.

```
<tag attribute="value" attribute2="value">Text content.</tag>
```

- Example page:

```
<html>
  <head>
    <title>My Web Page</title>
  </head>
  <body>
    <h1>Welcome To My Web Page</h1>
    
    <a href="pagelink.html" id="pageid">Check this other page.</a>
  </body>
</html>
```

In developing countries such as Brazil the deaths of impoverished infants are regularly unrecorded into the countries vital registration system; this causes a skew statistically. Cultural validity and contextual soundness can be used to ground the meaning of mortality from a statistical standpoint. In northeast Brazil they have accomplished this standpoint while conducting an ethnographic study combined with an alternative method to survey infant mortality.<sup>[62]</sup> These types of techniques can develop quality ethnographic data that will ultimately lead to a better portrayal of the magnitude of infant mortality in the region. Political economic reasons have been seen to skew the infant mortality data in the past when governor Ceans devised his presidency campaign on reducing the infant mortality rate during his term in office. By using this new way of surveying, these instances can be minimized and removed, overall creating accurate and sound data.<sup>[63]</sup>

However, IMR was, and remains, higher in LDCs. In 2001, the IMR for LDCs (91) was about 10 times as large as it was for MDCs (8). On average, for LDCs, the IMR is 17 times as higher than that of MDCs. Also, while both LDCs and MDCs made significant reductions in infant mortality rates, reductions among less developed

According to Guillot, Gerland, Pellerin and Saabneh 'birth histories, however, are subject to a number of errors, including omission of deaths and age misreporting errors.'<sup>100</sup>

Coronary artery research findings have demonstrated that nationwide racial differences in infant mortality are linked to the experiential status of the mother and that these disparities cannot be solely attributed to by socio-economic, behavioral or genetic factors.<sup>10,11</sup> The Hispanic population, an often observed in other health indicators, appears in the infant mortality rate, as well. Hispanic mothers are on average 30 percent less likely to be non-Hispanic white mothers, despite lower educational attainment and economic status. A study in North Carolina, for example, concluded that 'White women who did not graduate from high school have a lower infant mortality rate than black college graduates.<sup>12</sup> According to Mastaglio's (2004) Coronary Artery Risk Development in Young Adults' study, 'self-reported experiences of racial discrimination were associated with pre-term and low-birthweight deliveries, and these experiences may contribute to black-white disparities in prenatal outcomes.'<sup>13</sup> Likewise, dozens of population-based studies indicate that 'the subjective, or perceived experience of racial discrimination is strongly associated with an increased risk of infant death and with poor health prospects for future generations of African Americans.<sup>14</sup>

2050)				
Year	Rate	Year	Rate	
1950-1955	152	2000-2005	52	
1955-1960	138	2005-2010	47	
1960-1965	116	2010-2015	43	
1965-1970	110	2015-2020	40	
1970-1975	91	2020-2025	37	
1975-1980	83	2025-2030	34	
1980-1985	74	2030-2035	31	
1985-1990	66	2035-2040	28	
1990-1995	61	2040-2045	25	
1995-2000	57	2045-2050	23	



# Wikipedia on infant mortality

W Infant mortality - Wikipedia x +

← → ↻ https://en.wikipedia.org/wiki/Infant\_mortality

## Epidemiology [edit]

See also: *List of countries by infant mortality rate*

For the world, and for both less developed countries (LDCs) and more developed countries (MDCs), IMR declined. However, IMR was, and remains, higher in LDCs. In 2001, the IMR for LDCs (91) was about 10 times as large as that for MDCs. However, on average, much less than those among the more developed countries. <sup>[*clarification needed*]</sup>

A factor of about 67 separates the countries with the highest and lowest reported infant mortality rates. The top and

Rank	Country	Infant mortality rate (deaths/1,000 live births)
1	<a href="#">Afghanistan</a>	121.63
2	<a href="#">Niger</a>	109.98
3	<a href="#">Mali</a>	109.08
4	<a href="#">Somalia</a>	103.72
5	<a href="#">Central African Republic</a>	97.17
218	<a href="#">Sweden</a>	2.74
219	<a href="#">Singapore</a>	2.65
220	<a href="#">Bermuda</a>	2.47
221	<a href="#">Japan</a>	2.21
222	<a href="#">Monaco</a>	1.80

According to Guillot, Gerland, Pelletier and Saabneh "birth histories, however, are subject to a number of error

# Fetching a table from Wikipedia

```
library(rvest)

# 1) fetch and parse the website
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# 2) extract the html node containing the table
table <- html_node(page,
                    xpath = "//*[@id='mw-content-text']/div/table[2]")
# 3) extract the table as a data frame
mrates <- html_table(table)
```

# Inspecting the HTML source

- Convenient with modern browsers: Use the developer tools (right-click *Inspect*).
- Look at the HTML source to grasp the structure.
- Find out how to navigate the site.
- Find the element(s) you want to extract.
- Get the Xpath expression or CSS selector to extract elements.

## HTML elements visualized

w Infant mortality - Wikipedia x

← → ↻ 🔒 [https://en.wikipedia.org/wiki/Infant\\_mortality#Epidemiology](https://en.wikipedia.org/wiki/Infant_mortality#Epidemiology)

## Epidemiology [edit]

See also: *List of countries by infant mortality rate*

For the world, and for both less developed countries (LDCs) and more developed countries (MDCs), IMR declined significantly between 1960 and 2001. According to the [State of the World's Mothers report](#) by [Save the Children](#), the world IMR declined from 126 in 1960 to 57 in 2001.<sup>[10]</sup>

However, IMR was, and remains, higher in LDCs. In 2001, the IMR for LDCs (91) was about 10 times as large as it was for MDCs (8). On average, for LDCs, the IMR is 17 times as higher than that of MDCs. Also, while both LDCs and MDCs made significant reductions in infant mortality rates, reductions among less developed countries are, on average, much less than those among the more developed countries. <sup>[clarification needed]</sup>

A factor of about 67 separate countries with the highest and lowest reported infant mortality rates.

The top and bottom five countries by this measure (taken from The World Factbook's 2012

table.wikitable | 368 x 320  
be viewed here and shown below.

Rank	Country	Infant mortality rate (deaths/1,000 live births)
1	Afghanistan	121.63
2	Niger	109.98
3	Mali	109.08
4	Somalia	103.72
5	Central African Republic	97.17
216	Sweden	2.74
219	Singapore	2.65
220	Bermuda	2.47
221	Japan	2.21
222	Monaco	1.80

According to Guillot, Gerland, Pelletier and Saabneh "birth histories, however, are subject to a number of errors, including omission of deaths and age misreporting errors."<sup>103</sup>

United States [\[ edit \]](#)

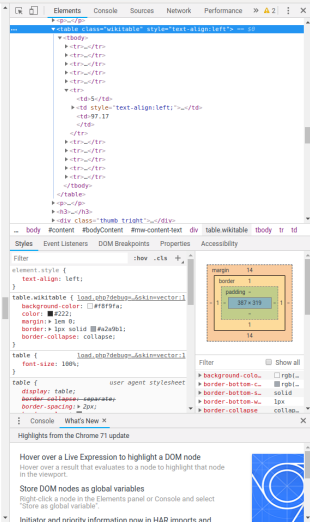
The infant mortality rate in the US decreased by 2.3% to a historic low of 582 infant deaths per 100,000 live births in 2014.<sup>[104]</sup>

Of the 27 most developed countries, the US has the highest Infant Mortality Rate, despite spending much more on health care per capita ([data needed](#)). Significant racial and socio-economic differences in the United States affect the IMR, in contrast with other developed countries, which have more homogeneous populations. In particular, IMR varies greatly by race in the US. The average IMR for the whole country is therefore not a fair representation of the wide variations that exist between segments of the population. Many theories have been explored as to why these racial differences exist, with socio-economic factors, such as income and education, responsible for a proportion. However, more studies

World historical and predicted infant mortality rates per 1,000 births (1950-2050)

UN, medium variant, 2008 rev.<sup>[100]</sup>

Years	Rate	Years	Rate
1950–1955	152	2000–2005	52
1955–1960	136	2005–2010	47
1960–1965	116	2010–2015	43
1965–1970	100	2015–2020	40
1970–1975	91	2020–2025	37
1975–1980	83	2025–2030	34
1980–1985	74	2030–2035	31
1985–1990	65	2035–2040	28
1990–1995	61	2040–2045	25
1995–2000	57	2045–2050	23



# Infant mortality rates from Wikipedia

```
<table class="wikitable" style="text-align:left">
  <tbody>
    <tr>
      <th>Rank</th>
      <th>Country</th>
      <th>Infant mortality rate <br> (deaths/1,000 live births)</th>
    </tr>
    <tr>
      <td>1</td>
      <td style="text-align:left;"><a href="/wiki/Afghanistan" title="Afghanistan">Afghanistan</a></td>
      <td>121.63</td>
    </tr>
    <tr>
      <td>2</td>
      <td style="text-align:left;"><a href="/wiki/Niger" title="Niger">Niger</a></td>
      <td>109.98</td>
    </tr>
    <tr>
      <td>3</td>
      <td style="text-align:left;"><a href="/wiki/Mali" title="Mali">Mali</a></td>
      <td>109.08</td>
    </tr>
    <tr>
      <td>4</td>
      <td style="text-align:left;"><a href="/wiki/Somalia" title="Somalia">Somalia</a></td>
      <td>103.72</td>
    </tr>
    ...
  </tbody>
</table>
```



# CSS selectors and XPath expressions

```
# fetch and parse the website
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# list the table nodes
html_nodes(page, "table")
# using xpath expressions or css selectors is equivalent
table <- html_node(page,
                    xpath = "//*[@id='mw-content-text']/div/table[2]")
table <- html_node(page,
                    css = "#mw-content-text > div > table:nth-child(121)")
```

## Examples

---

# The general structure

- There is no universal recipe. But most programs follow a certain structure.
  1. Open a website mimicking a browser and navigate it (optional).
  2. Get the page source HTML and feed it to a parser.
  3. Extract the elements you need.
  4. Filter and arrange them as needed and save them.
  5. Repeat 1.–4. as needed.

# Navigating to another page

```
# Open infant mortality page
session <- html_session("https://en.wikipedia.org/wiki/Infant_mortality")
# Goto page on Somalia
session <- follow_link(session, "Somalia")
# Read the source
page <- read_html(session)
# Extract html
table <- html_node(page,
  xpath = "//*[@id='mw-content-text']/div/table[4]")
regions <- html_table(table)
```

# Filtering links

```
# read wiki page
page <- read_html("https://en.wikipedia.org/wiki/Infant_mortality")
# get the links
wikilinks <- html_attr(html_nodes(page, "a"), "href")
# use regex to filter internal links:
#   select only articles, no files or category pages,
#   matching with mortality or somalia
links <- grep("^(?!.*:)(/wiki/. *Mortality)|(/wiki/. *Somalia)", wikilinks,
              ignore.case = TRUE, value = TRUE, perl = TRUE)
links <- unique(links)
# go to first selected article page and process it
session <- jump_to(session, links[1])
page <- read_html(session)
html_nodes(page, "title")
```

## A more elaborate example

- Phillippine Statistics Authority *Good Governance Index*.
- Available at <http://nap.psa.gov.ph/ggi/default.asp>.
- Available at <https://web.archive.org/web/20190915135458/http://nap.psa.gov.ph/ggi/default.asp>.
- Get all GGI data tables for all municipalities.
- Save them in a local data file for further analysis.
- How would you go about this?

## If simple navigation fails

- Some web pages cannot be navigated easily with simpler requests.
- Often due to hidden Javascript or other server-side processing.
- In this case, resort to Selenium (`library(rselenium)` in R).
- Under the hood, Selenium relies on a complete browser running in a container.
- Slower and comes with substantial overhead costs.
- Only use when absolutely necessary.

## Another example

- WHO venomous snakes distribution and species risk categories
- Database form link:  
<https://apps.who.int/bloodproducts/snakeantivenoms/database/SearchFrm.aspx>
- Collect snake data for all countries.
- Getting dropdown options and initial form submission straightforward.
- So is table extraction.
- Navigating further links is tricky.



## General remarks

- Start simple and expand your program incrementally.
- Keep it simple. Do not overengineer the problem. Do not repeat yourself.
- Limit the number of iterations for test runs. Use print statements to inspect objects.
- Write tests to verify things work as intended.
- If your program requires complex monitoring/validation procedures or threading for performance, use Python.

## Final remarks

- Sometimes small programs can go a long way.
- Do not lose sight of your ultimate goal. Time is valuable.
- Do not engage in perfectionism, focus on GTD.
- Identify everyday tasks that you can optimize.
- It might even be fun.

# Assignment

---

# Assignment, part I

Choose *one* of the following.

## Open Data Tanzania

<https://tanzania.opendataforafrica.org/>

- ➡ For each region in Tanzania, get the basic facts (top of the page: capital, population, area, ...) and mortality statistics (table under link 'Deaths').

## WHO Venomous Snakes Database

<http://apps.who.int/bloodproducts/snakeantivenoms/database/SearchFrm.aspx>

- ➡ For every country, collect *all* species of venomous snakes.

## Swiss Tax Calculator

<https://swisstaxcalculator.estv.admin.ch/#/taxburden/income-wealth-tax>

- ➡ Collect municipal income tax burden statistics for the following parameters: All payers, atheist, geographical comparison (all municipalities), all years, income 50'000 CHF.

# Assignment, part I

- Submission deadline is October 25, 2020.
- Submit code only, no data.
- Comment your code or submit a short description alongside.
- A proof-of-concept restricted to the first couple of regions/... is fine.
- Accounts for 20% of the final grade.

Next lecture: Text as data.