

Week 1: Introduction to QTA

POP77032 Quantitative Text Analysis for Social Scientists

Tom Paskhalis

Overview

- Text as data
- Data collection
- Web technologies
- HTML fundamentals
- XPath

Text As Data

Textual Data

- Ubiquitous
- Yet often underutilized



Elon Musk @elonmusk · 13 May

Twitter deal temporarily on hold pending details supporting calculation that spam/fake accounts do indeed represent less than 5% of users



Τα σόσιαλ μίντια καταστρέφουν τη Δημοκρατία

Η δημοσιογράφος, συγγραφέας και εδώ και λίγες μέρες κάτοχος του Νόμπελ Ειρήνης Μαρία Ρέσα μιλά στην Καθημερινή



Newsroom

24.10.2021 · 12:10

Η δημοσιογράφος, συγγραφέας και εδώ και λίγες μέρες κάτοχος του Νόμπελ Ειρήνης Μαρία Ρέσα μίλησε στην Καθημερινή για το πικρό μήνυμα που στέλνει η βράβευσή της και του Ρώσου δημοσιογράφου Ντμίτρι Μουράτοφ σε μια εποχή που η Δημοκρατία απειλείται από την άνοδο ακραίων δυνάμεων και τη διάδοση των fake news μέσω των σόσιαλ μίντια.

Web Scraping

Online Data Sources

- Data downloadable in tabular format (E.g. CSV/TSV, XLS, DTA, etc.)
- Data available online as a table (E.g. webpages with rendered tables)
- Unstructured data available online (E.g. simple webpages)
- Interactive webpages with user-input (E.g. webpages with logins, dropdown menus)
- Web APIs (special interfaces for querying, e.g. Twitter, Google)

Online Data Collection

- Tabular format: download single or multiple files (automate with `download.file()` in R, `wget` in Python/Terminal)
- Online tables and unstructured data: simple web scraping (HTML with XPath, `rvest` in R, `beautifulsoup` in Python)
- Interactive webpages: web scraping with headless browser (Selenium, Playwright - Python bindings recommended)
- Web API: sending requests and processing responses (HTTP queries, `httr2` in R, `requests` in Python)

(Wikipedia)

8

Unstructured Data

https://eur-lex.europa.eu/search.html?DTS_SUBDOM=ALL_ALL&DTS_DOM=ALL&type=advanced&excConsLeg=true&qid=1638131808500&SUBDOM_INIT=ALL_ALL&page=

An official website of the European Union How do you know?

EUR-Lex
Access to European Union law

English My EUR-Lex
Experimental features

MENU QUICK SEARCH

Search tips Need more search options? Use the Advanced search

EUROPA > EUR-Lex home > Advanced search > Search results

Search Results

You can only view pages 1–9,999 of the search results.

Search criteria

Results 1 - 10 of 1021911 Sort by Relevant

1 2 > >>

Clear selection Customise shown information Export

☐ Directive (EU) 2021/1883 of the European Parliament and of the Council of 20 October 2021 on the conditions of entry and residence of third-country nationals for the purpose of highly qualified employment, and repealing Council Directive 2009/50/EC
PE/40/2021/REV/1
OJ L 382, 28.10.2021, p. 1–38 (BG, ES, CS, DA, DE, ET, EL, EN, FR, GA, HR, IT, LV, LT, HU, MT, NL, PL, PT, RO, SK, SL, FI, SV)
In force
CELEX number: 32021L1883
Author: European Parliament, Council of the European Union
Date of document: 20/10/2021; Date of signature

☐ Regulation (EU) 2021/1873 of the European Parliament and of the Council of 20 October 2021 on the extension of the term of the Community plant variety rights for varieties of the species *Asparagus officinalis* L. and of the species groups flower bulbs, woody small fruits and woody ornamentals
PE/50/2021/REV/2
OJ L 378, 26.10.2021, p. 1–3 (BG, ES, CS, DA, DE, ET, EL, EN, FR, GA, HR, IT, LV, LT, HU, MT, NL, PL, PT, RO, SK, SL, FI, SV)
In force
CELEX number: 32021R1873
Author: European Parliament, Council of the European Union
Date of document: 20/10/2021; Date of signature

Refine query

You have selected:

All

By keyword

In title In text

By year of document

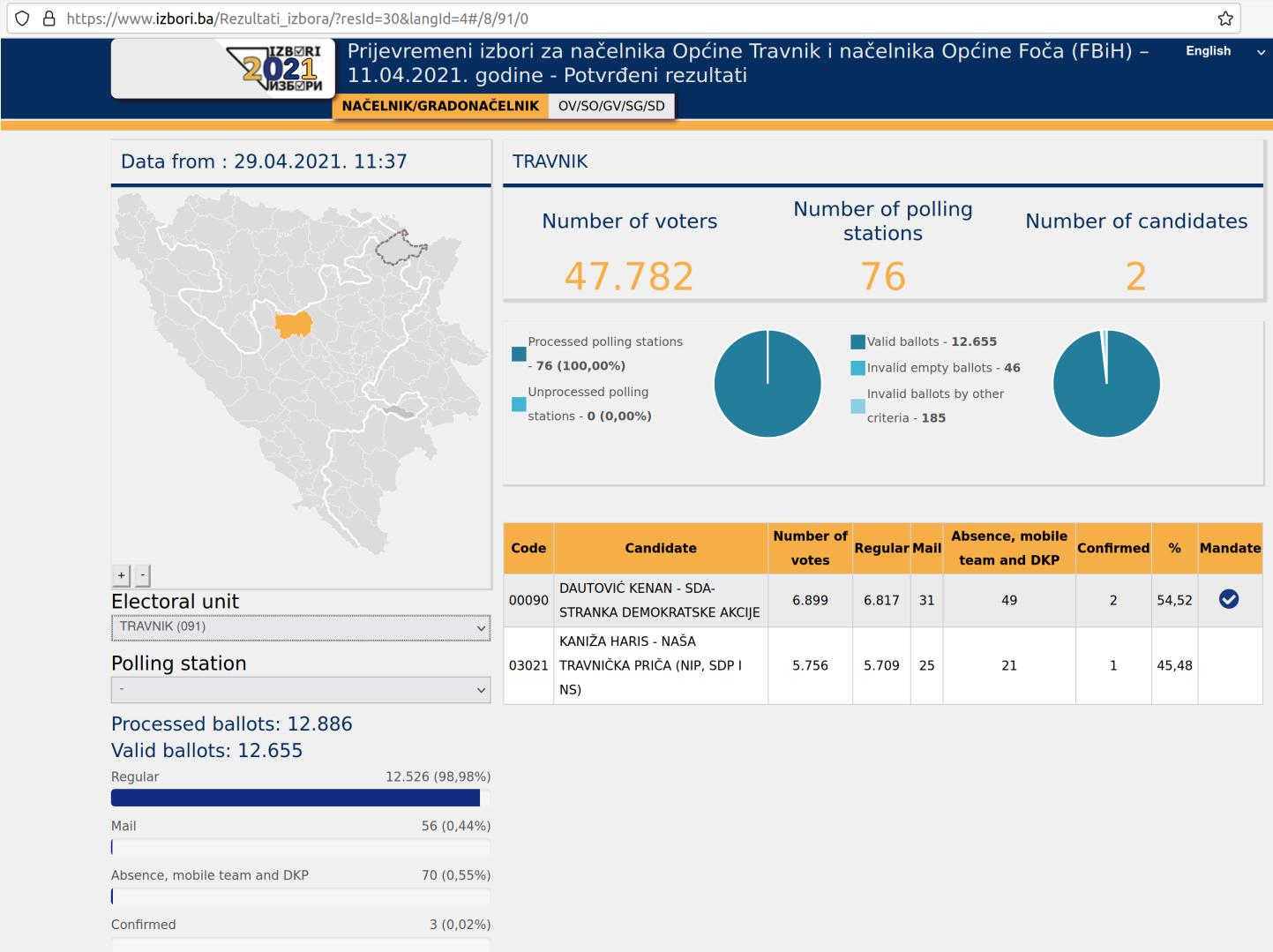
2021 (17116)
2020 (20801)
2019 (21927)
2018 (23944)
2017 (20997)
See more...

By Collection

EU law and case-law (672958)
Legal acts (218081)
Treaties (9810)
International agreements (11816)
Preparatory documents (133709)
Parliamentary questions (197036)
Case-law (100424)
EFTA documents (2082)
National law and case-law (215568)
National transposition (177366)

(Eur-Lex)

Interactive Webpages



(Izbori.ba)

Automated Data Collection

- Manual scraping (copy-pasting) can be:
 - Extremely laborious and time-consuming
 - Very error-prone
 - Often impossible to reproduce exactly
- Automated data collection
 - Easy to scale up (computer time is cheap)
 - Less error-prone
 - Usually, perfectly reproducible
- There is a trade-off (time invested in automation vs time saved)
 - However, it is good to err on the side of automation

Web Technologies

- Key technologies used to disseminate content on the Web:
 - XML/HTML (Extensible Markup Language/Hypertext Markup Language)
 - CSS (Cascading Style Sheets)
 - JavaScript
 - API (Application Programming Interface)
 - JSON (JavaScript Object Notation)

Static vs Dynamic Websites

- The critical feature of a website which determines approach to scraping its content
- *Static* websites all have prebuild source code which is served at user's request
 - No real-time processing of user's input
 - They can contain elements that change the appearance of a website
 - Example: POP77142 website
- *Dynamic* websites render websites in real-time as a response to user's input
 - They can use a range of technologies to achieve it (JavaScript, Python Django, PHP)
 - Example: Google Maps

HTML: Hypertext Markup Language

- HTML (**H**ypertext **M**arkup **L**anguage) is a mark-up language for webpages
- Forms the basis of static websites
- Your browser renders (interprets) HTML for viewing
- Current version is HTML5

```
<!DOCTYPE html>
<html>
  <head>
    <title>A title</title>
  </head>
  <body>
    <h1 style="color:Red;">A heading</h1>
    <p>A paragraph.</p>
  </body>
</html>
```

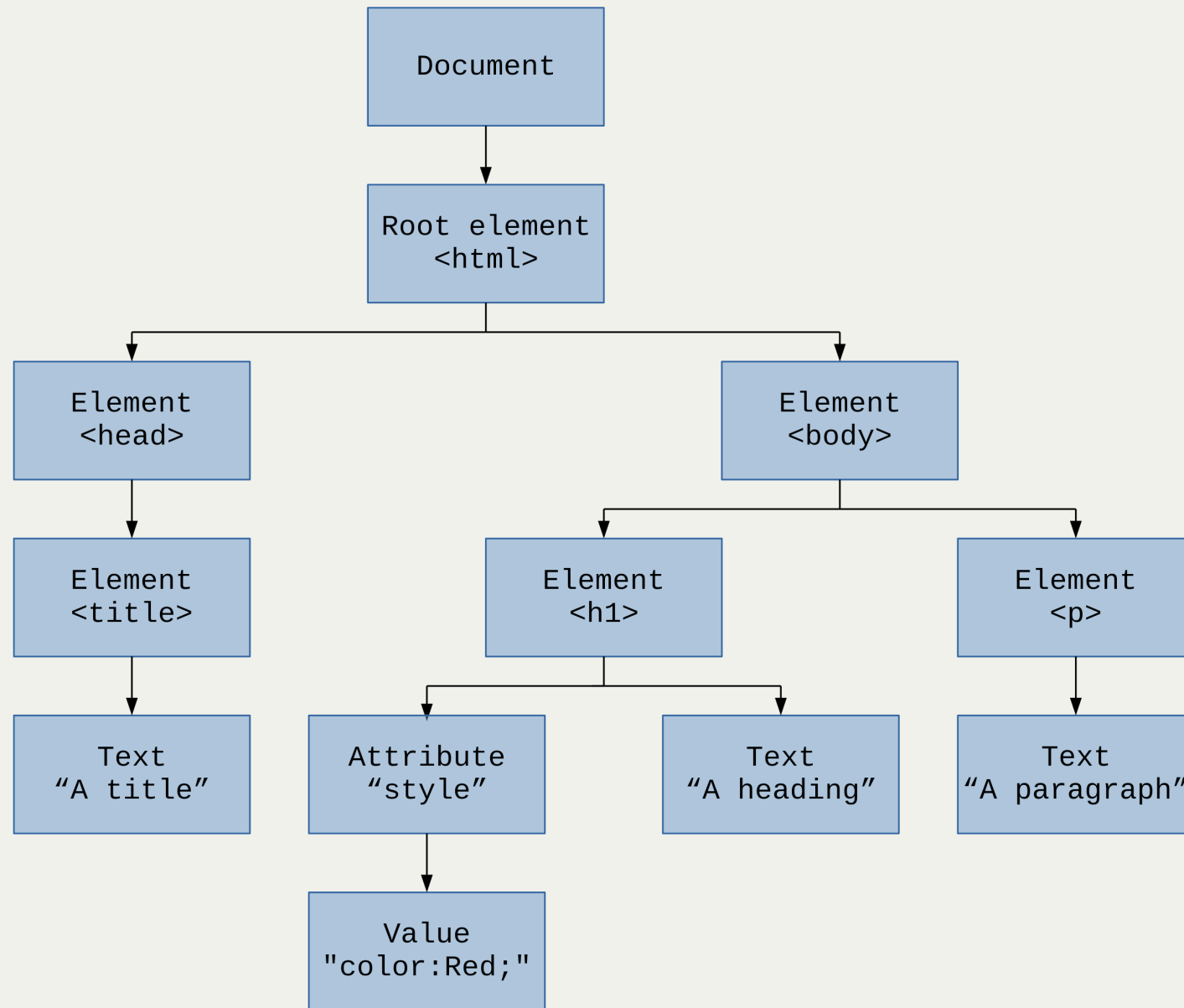


W3Schools: Try HTML

HTML Basics

- Basic unit of HTML is an *element* (aka *node*)
- Elements, typically, begin with an *start tag* (e.g. `<h1>`)
- And finish with an *end tag* (e.g. `</h1>`)
- Content of an element is found between the start and end tags
- *Attributes* are special words used within a start tag to control element's behaviour (e.g. `style="color:Red;"`)
- Some HTML tag examples:
 - Document structure: `<html>`, `<body>`, `<header>`
 - Document components: `<h1>`, `<title>`, `<div>`
 - Text style: ``, `<i>`
 - Hyperlinks: `<a>`

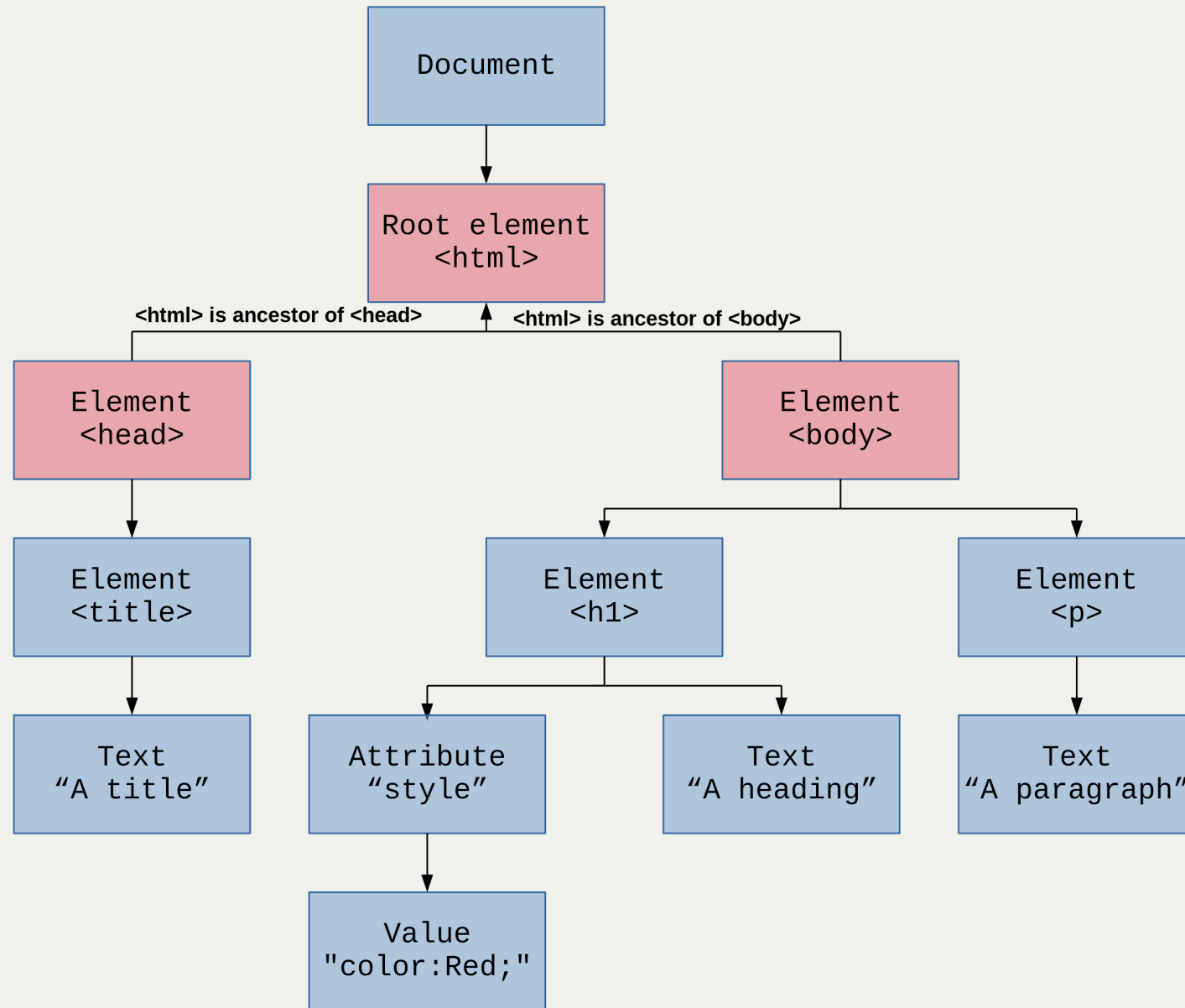
HTML tree



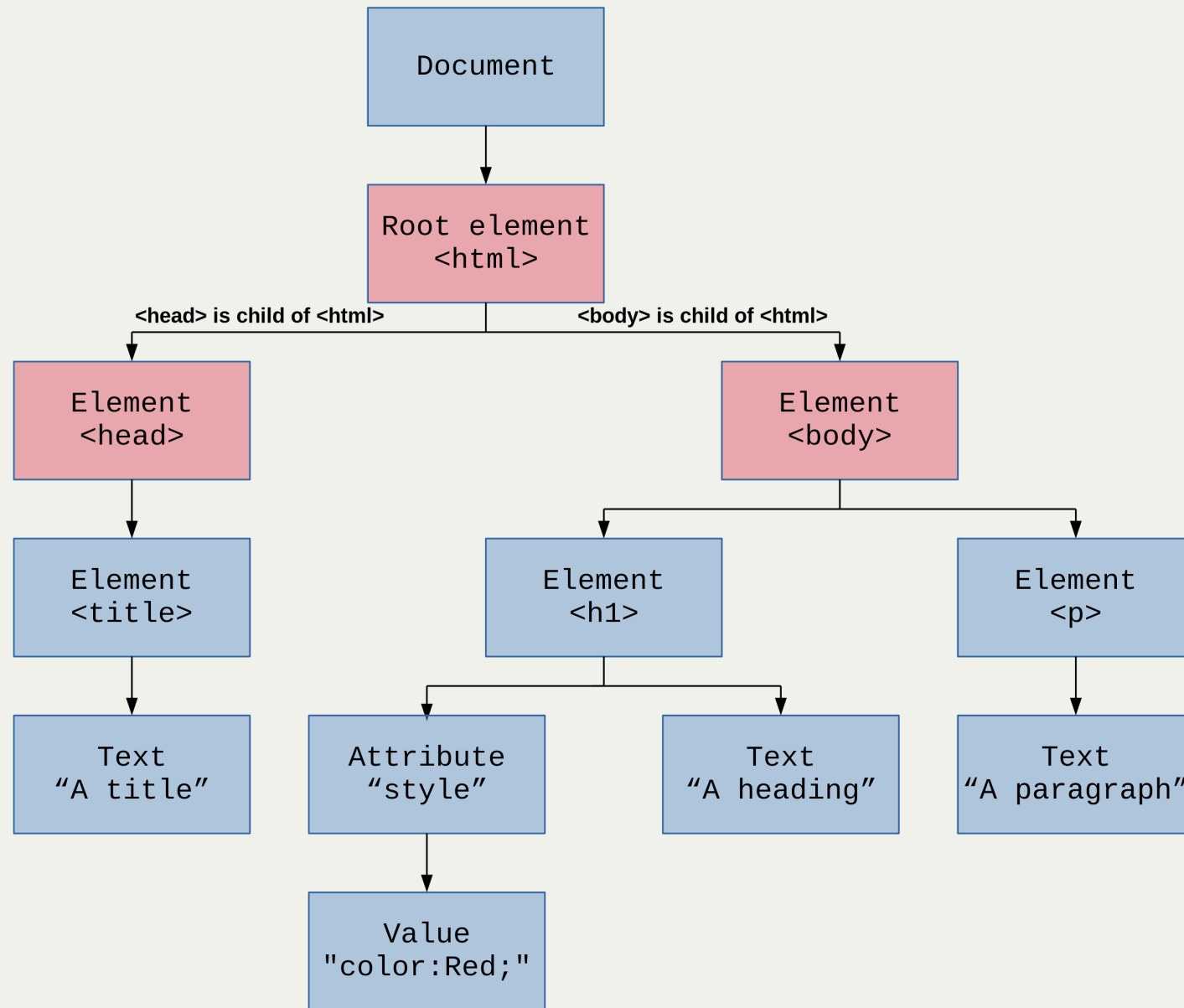
HTML Tree Relationships

- All elements (nodes) in HTML tree are connected by relationships
- These relationship can be of the following types:
 - Ancestors (parents)
 - Descendants (children)
 - Siblings

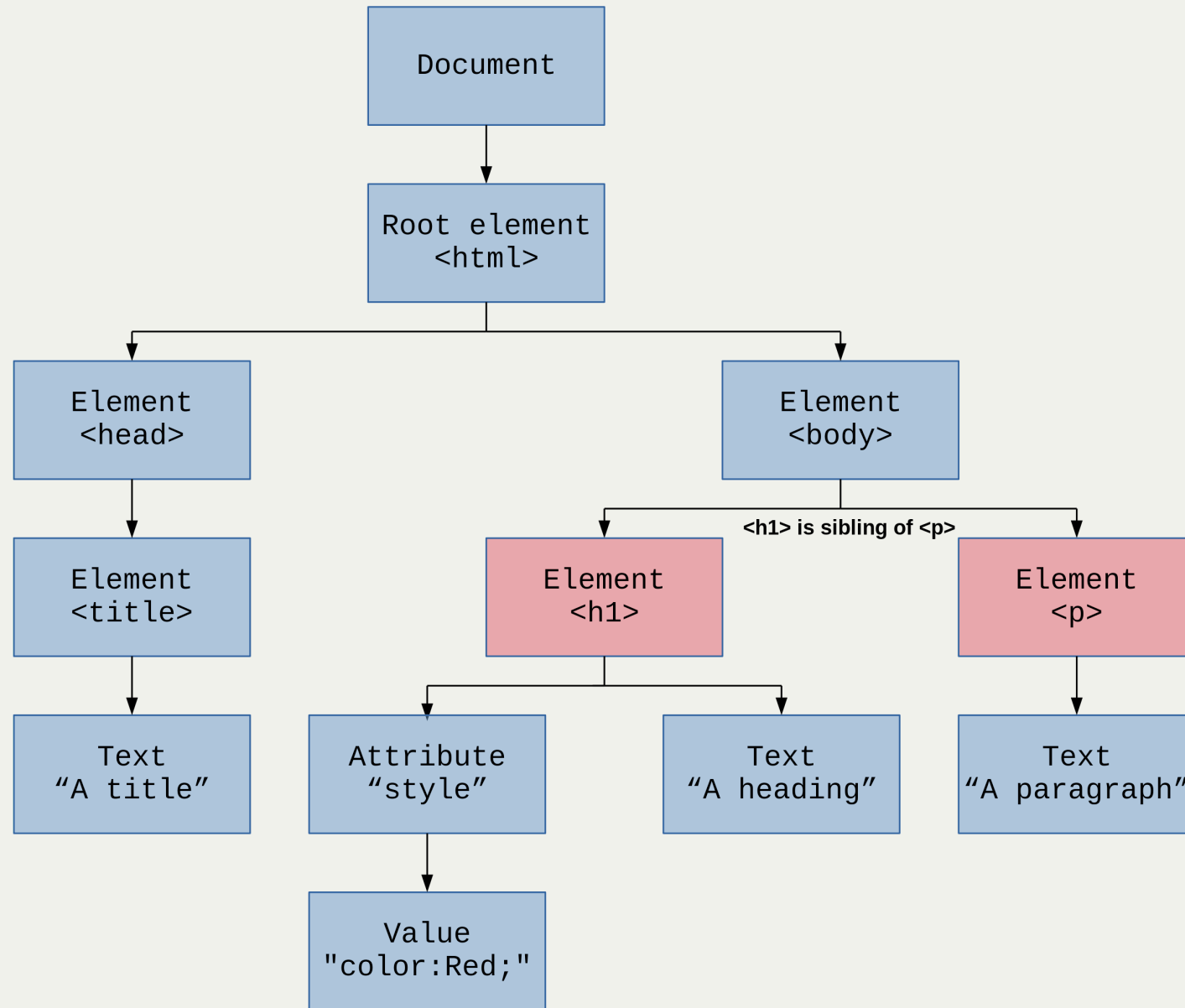
HTML Parent/Ancessor



HTML Children/Descendants



HTML Siblings



Parsing HTML Tree: Example

```
1 library("rvest")
```

```
1 html_txt <- "  
2 <!DOCTYPE html>  
3 <html>  
4   <head>  
5     <title>A title</title>  
6   </head>  
7   <body>  
8     <h1 style='color:Red;'>A heading</h1>  
9     <p>A paragraph.</p>  
10  </body>  
11 </html>"
```

```
1 html <- rvest::read_html(html_txt)
```

```
1 str(html)
```

List of 2

\$ node:<externalptr>

\$ doc :<externalptr>

- attr(*, "class")= chr [1:2] "xml_document" "xml_node"

Parsing HTML Tree: Example

```
1 children <- rvest::html_children(html)
2 children
```

```
{xml_nodeset (2)}
[1] <head>\n<meta http-equiv="Content-Type" content="text/html; charset=UTF-8 ...
[2] <body>\n      <h1 style="color:Red;">A heading</h1> \n      <p>A para ...
```

```
1 body <- children[2]
2 rvest::html_name(body)
```

```
[1] "body"
```

```
1 children2 <- rvest::html_children(body)
2 children2
```

```
{xml_nodeset (2)}
[1] <h1 style="color:Red;">A heading</h1>
[2] <p>A paragraph.</p>
```

```
1 rvest::html_attrs(children2[1])
```

```
[[1]]
      style
"color:Red;"
```

```
1 rvest::html_text(children2[1])
```

```
[1] "A heading"
```

XML: Extensible Markup Language

- XML (Extensible Markup Language) is a more general form of markup language
- Allows sharing structured data of tree-like form
- Relative to HTML:
 - Tags are user-defined
 - End tags are always required
 - Stricter (no inconsistencies permitted)

```
<?xml version="1.0" encoding="UTF-8" ?>
<courses>
  <course>
    <title>Computer Programming for Social Scientists</title>
    <code>POP77001</code>
    <year>2024</year>
    <term>Michaelmas</term>
    <description>Course on computer programming in Python and R.
  </description>
  </course>
  <course>
    <title>Quantitative Text Analysis for Social
Scientists</title>
    <code>POP77142</code>
    <year>2025</year>
    <term>Hillary</term>
    <description>Introduction to text analysis.</description>
  </course>
</courses>
```

Parsing XML Tree: Example

```
1 library("xml2")

1 xml_txt <-
2 '<?xml version="1.0" encoding="UTF-8" ?>
3 <courses>
4   <course>
5     <title>Computer Programming for Social Scientists</title>
6     <code>POP77001</code>
7     <year>2024</year>
8     <term>Michaelmas</term>
9     <description>Course on computer programming in Python and R.</description>
10  </course>
11  <course>
12    <title>Quantitative Text Analysis for Social Scientists</title>
13    <code>POP77142</code>
14    <year>2025</year>
15    <term>Hillary</term>
16    <description>Introduction to text analysis.</description>
17  </course>
18 </courses>'

1 xml <- xml2::read_xml(xml_txt)

1 str(xml)
```

List of 2

```
$ node:<externalptr>
$ doc :<externalptr>
- attr(*, "class")= chr [1:2] "xml_document" "xml_node"
```


Parsing XML Tree: Example

```
1 children3 <- xml2::xml_children(xml)
2 children3
```

```
{xml_nodeset (2)}
[1] <course>\n  <title>Computer Programming for Social Scientists</title>\n  ...
[2] <course>\n  <title>Quantitative Text Analysis for Social Scientists</titl ...
```

```
1 pop77001 <- children3[1]
2 xml2::xml_children(pop77001)
```

```
{xml_nodeset (5)}
[1] <title>Computer Programming for Social Scientists</title>
[2] <code>POP77001</code>
[3] <year>2024</year>
[4] <term>Michaelmas</term>
[5] <description>Course on computer programming in Python and R.</description>
```

```
1 xml2::xml_text(xml_children(children3[1]))
```

```
[1] "Computer Programming for Social Scientists"
[2] "POP77001"
[3] "2024"
[4] "Michaelmas"
[5] "Course on computer programming in Python and R."
```

Examples of XML

- RSS (**R**eally **S**imple **S**yndication) feeds
- SVG (**S**calable **V**ector **G**raphics) images
- Modern office documents (Microsoft Office `.docx`, `.xlsx`, `.pptx`, OpenOffice/LibreOffice)

Parsing XML/HTML with XPath

- XPath (XML Path Language) is a language for selecting parts of XML/HTML tree
- Basic syntax:
 - `/` - select element at the root node (e.g. `/html/body`)
 - `//` - select element at any depth (e.g. `//h1`)
 - `//<tag>/*` - select all descendants of tag (e.g. `//body/*`)
 - `//<tag>[@<attr>]` - select all elements that have given attribute (e.g. `//h1[@style]`)
 - `//<tag>[@<attr>='<value>']` - select all elements, whose attribute has given value (e.g. `//h1[@style='color:Red;']`)



Extra

XPath syntax

Parsing XML/HTML with XPath: Example

```
1 rvest::html_elements(html, xpath = "//p")
```

```
{xml_nodeset (1)}  
[1] <p>A paragraph.</p>
```

```
1 rvest::html_elements(html, xpath = "//h1[@style='color:Red;']")
```

```
{xml_nodeset (1)}  
[1] <h1 style="color:Red;">A heading</h1>
```

```
1 xml2::xml_find_all(xml, xpath = "//code")
```

```
{xml_nodeset (2)}  
[1] <code>POP77001</code>  
[2] <code>POP77142</code>
```

```
1 # We can also find elements by text  
2 xml2::xml_find_all(xml, xpath = "//code[text()='POP77001']")
```

```
{xml_nodeset (1)}  
[1] <code>POP77001</code>
```

Scraping Webpage

←

→

↻

en.wikipedia.org/wiki/Members_of_the_1st_Dáil

📄

80%

☆

📄

🔍

WIKIPEDIA

25 years of the free encyclopedia

🔍

Search Wikipedia

Search

Members of the 1st Dáil

2 languages

Contents

hide

(Top)

Composition of the 1st Dáil

Members by constituency

Changes

Vacancies

By-elections

See also

Notes

References

Members of the 1st Dáil

Talk

Read

Edit

View history

Tools

From Wikipedia, the free encyclopedia

This article **needs additional citations for verification**. Please help [improve this article](#) by [adding citations to reliable sources](#). Unsourced material may be challenged and removed.

Find sources: "Members of the 1st Dáil" – news · newspapers · books · scholar · JSTOR (May 2014) [\(Learn how and when to remove this message\)](#)

The members of the **First Dáil**, known as *Teachtai Dála* (TDs), were the 101^[a] **Members of Parliament** (MPs) **returned from constituencies in Ireland** at the **1918 United Kingdom general election**. In its first general election, Sinn Féin won 73^[a] seats and viewed the result as a mandate for independence; in accordance with its declared policy of **abstentionism**, its 69^[a] MPs refused to attend the **British House of Commons** in Westminster, and established a revolutionary parliament known as **Dáil Éireann**. The other Irish MPs — 26 **unionists** and six^[b] from the **Irish Parliamentary Party** (IPP) — sat at Westminster and for the most part ignored the invitation to attend the Dáil. **Thomas Harbison**, IPP MP for **North East Tyrone**, did acknowledge the invitation, but "stated he should decline for obvious reasons".^[1] The Dáil met for the first time on 21 January 1919 in **Mansion House** in **Dublin**. Only 27 members attended; most of the other Sinn Féin TDs were imprisoned by the British authorities, or in hiding under threat of arrest. All 101 MPs were considered TDs, and their names were called out on the roll of membership, though there was some laughter when **Irish Unionist Alliance** leader **Edward Carson** was described as *as láthair* ("absent").^[2] The database of members of the **Oireachtas** (Irish parliament) includes for the First Dáil only those elected for Sinn Féin.^[3]

Composition of the 1st Dáil

edit

Party	Dec. 1918 ^[c]	May 1921 ^[d]
<div><div></div></div> Sinn Féin	73 ^[a]	69
<div><div></div></div> Irish Unionist	22	—
<div><div></div></div> Irish Parliamentary	6 ^[b]	2
<div><div></div></div> Labour Unionist	3	—
<div><div></div></div> Ind. Unionist	1	1
<div><div></div></div> UUP	—	23
<div><div></div></div> Nationalist	—	4
<div><div></div></div> Unionist Anti-Partition League	—	2
<div><div></div></div> Vacant	—	4
Total		105

Government party denoted with bullet (●).

Members by constituency

edit

Constituency	Name	Portrait	Party affiliation		Assumed office
			Start of Dáil term	End of Dáil term	
Antrim East	Robert McCalmont	<div><div></div></div>	<div><div></div></div> Irish Unionist	Resigned in 1919	Abstained
	George Hanna	<div><div></div></div>	Elected in 1919 by-election as Independent Unionist	<div><div></div></div> Ulster Unionist	Abstained

1st Dáil

2nd Dáil

24 of the 27 TDs present at the first Dáil meeting on 21 January 1919, photographed afterwards on the steps of the Mansion House. The caption gives names in Irish.

Overview

Legislative body

Dáil Éireann

Jurisdiction

Irish Republic

Meeting place

Mansion House

UCD (Earlsfort Terrace)

Term

21 January 1919 – 10 May 1921

Election

1918 general election

Government

1st Dáil ministry

(until 22 January 1919)

2nd Dáil ministry

(1919–1922)

Members

105^[a]

Ceann Comhairle

Seán T. O'Kelly

— Count Plunkett

22 January 1919

— Cathal Brugha

until 22 January 1919

President of Dáil Éireann

Éamon de Valera

— Cathal Brugha

until 1 April 1919

(Wikipedia)

29

Scraping Webpage with XPath: Example

```
1 html <- rvest::read_html("https://en.wikipedia.org/wiki/Members_of_the_1st_D%C3%A1il")
```

```
1 tables <- rvest::html_elements(html, xpath = "//table")
2 tables
```

```
{xml_nodeset (8)}
[1] <table class="box-More_citations_needed plainlinks metadata ambox ambox-c ...
[2] <table class="infobox vevent"><tbody>\n<tr><th colspan="2" class="infobox ...
[3] <table style="width:100%; border-collapse:collapse"><tbody><tr style="ver ...
[4] <table class="wikitable" style="font-size: 95%;"><tbody>\n<tr>\n<th colsp ...
[5] <table class="wikitable sortable"><tbody>\n<tr>\n<th rowspan="2">Constitu ...
[6] <table class="wikitable"><tbody>\n<tr>\n<th>Constituency\n</th>\n<th>Outg ...
[7] <table class="wikitable"><tbody>\n<tr>\n<th>Winner\n</th>\n<th colspan="2 ...
[8] <table class="nowraplinks mw-collapsible autocollapse navbox-inner" style ...
```

```
1 tbody <- rvest::html_children(tables[5])
2 tbody
```

```
{xml_nodeset (1)}
[1] <tbody>\n<tr>\n<th rowspan="2">Constituency\n</th>\n<th rowspan="2">Name\ ...
```

```
1 tds <- rvest::html_table(tbody)
2 tds
```

```
[[1]]
# A tibble: 109 × 8
  Constituency Name Portrait `Party affiliation` `Party affiliation`
  <chr> <chr> <chr> <chr> <chr>
1 Constituency Name "Portra... "Start of Dáil ter... Start of Dáil term
2 Antrim East Robert McCal... "" "" Irish Unionist
```

3	Antrim East	George Hanna	""	"Elected in 1919 b...	Elected in 1919 by...
4	Antrim Mid	Hugh O'Neill	""	""	Irish Unionist
5	Antrim North	Peter Kerr-S...	""	""	Irish Unionist
6	Antrim South	Charles Craig	""	""	Irish Unionist
7	Armagh Mid	James Lonsda...	""	""	Irish Unionist
8	Armagh North	William Allen	""	""	Irish Unionist
9	Armagh South	Patrick Donn...	""	""	Irish Parliamentary
10	Belfast Cromac	William Arth...	""	""	Irish Unionist

i 99 more rows

i 3 more variables: `Party affiliation` <chr>, `Party affiliation` <chr>,
`Assumed office` <chr>

Scraping Webpage with XPath:

Example

```
1 str(tds)
```

List of 1

```
$ : tibble [109 × 8] (S3: tbl_df/tbl/data.frame)
  ..$ Constituency      : chr [1:109] "Constituency" "Antrim East" "Antrim East" "Antrim Mid" ...
  ..$ Name              : chr [1:109] "Name" "Robert McCalmont" "George Hanna" "Hugh O'Neill" ...
  ..$ Portrait          : chr [1:109] "Portrait" "" "" "" ...
  ..$ Party affiliation: chr [1:109] "Start of Dáil term" "" "Elected in 1919 by-electionas Independent
Unionist" "" ...
  ..$ Party affiliation: chr [1:109] "Start of Dáil term" "Irish Unionist" "Elected in 1919 by-electionas
Independent Unionist" "Irish Unionist" ...
  ..$ Party affiliation: chr [1:109] "End of Dáil term" "Resigned in 1919" "" "" ...
  ..$ Party affiliation: chr [1:109] "End of Dáil term" "Resigned in 1919" "Ulster Unionist" "Ulster Unionist"
...
  ..$ Assumed office    : chr [1:109] "Assumed office" "Abstained" "Abstained" "Abstained" ...
```

```
1 tds <- tds[[1]]
2 head(tds)
```

A tibble: 6 × 8

Constituency	Name	Portrait	`Party affiliation`	`Party affiliation`
<chr>	<chr>	<chr>	<chr>	<chr>
1 Constituency	Name	"Portra...	"Start of Dáil ter...	Start of Dáil term
2 Antrim East	Robert McCalmont	""	""	Irish Unionist
3 Antrim East	George Hanna	""	"Elected in 1919 b...	Elected in 1919 by...
4 Antrim Mid	Hugh O'Neill	""	""	Irish Unionist
5 Antrim North	Peter Kerr-Smil...	""	""	Irish Unionist
6 Antrim South	Charles Craig	""	""	Irish Unionist

Scraping Webpage with XPath:

Example

```
1 colnames(tds) <- tds[1,]
2 tds <- tds[-1,]
3 head(tds)
```

A tibble: 6 × 8

	Constituency	Name	Portrait	`Start of Dáil term`	`Start of Dáil term`
	<chr>	<chr>	<chr>	<chr>	<chr>
1	Antrim East	Robert McCalm...	""	""	Irish Unionist
2	Antrim East	George Hanna	""	"Elected in 1919 by...	Elected in 1919 by-...
3	Antrim Mid	Hugh O'Neill	""	""	Irish Unionist
4	Antrim North	Peter Kerr-Sm...	""	""	Irish Unionist
5	Antrim South	Charles Craig	""	""	Irish Unionist
6	Armagh Mid	James Lonsdale	""	""	Irish Unionist

i 3 more variables: `End of Dáil term` <chr>, `End of Dáil term` <chr>,
`Assumed office` <chr>

```
1 tds <- tds[, -3]
2 str(tds)
```

tibble [108 × 7] (S3: tbl_df/tbl/data.frame)

```
$ Constituency      : chr [1:108] "Antrim East" "Antrim East" "Antrim Mid" "Antrim North" ...
$ Name              : chr [1:108] "Robert McCalmont" "George Hanna" "Hugh O'Neill" "Peter Kerr-Smiley" ...
$ Start of Dáil term: chr [1:108] "" "Elected in 1919 by-electionas Independent Unionist" "" "" ...
$ Start of Dáil term: chr [1:108] "Irish Unionist" "Elected in 1919 by-electionas Independent Unionist" "Irish
Unionist" "Irish Unionist" ...
$ End of Dáil term  : chr [1:108] "Resigned in 1919" "" "" "" ...
$ End of Dáil term  : chr [1:108] "Resigned in 1919" "Ulster Unionist" "Ulster Unionist" "Ulster Unionist" ...
$ Assumed office    : chr [1:108] "Abstained" "Abstained" "Abstained" "Abstained" ...
```

Web Scraping in Practice

- Always check first whether an API for querying exists.
- It is the most robust (and sanctioned) way of obtaining data.
- Check copyrights and respect those when using scraped data.
- Limit you scraping bandwidth (introduce waiting times between queries).

Next

- Tutorial: HTML and web scraping