

# Web Scraping mit R

Stefan Schmidt

16.04.2020

# Agenda

- ▶ Web Scraping
- ▶ Review: HTML<sup>1</sup> und CSS<sup>2</sup> Selektoren
- ▶ R Package: rvest
- ▶ Werkzeug: selectorGadget
- ▶ Rechtliche Bedenken

---

<sup>1</sup>Hypertext Markup Language

<sup>2</sup>Cascading Style Sheets

# Web Scraping

# Web Scraping

Was ist das?

- ▶ engl. “aus dem Web kratzen” (auch “Web **Harvesting**”)
- ▶ Verfahren zum automatisierten Auslesen von Texten / Daten von Webseiten (**HTML**)

# Web Scraping

Wann wendet man es an?

- ▶ Wenn zur Extraktion der Informationen keine API<sup>3</sup> zur Verfügung steht
- ▶ Also keine Möglichkeit besteht definiert und sauber (z.B. via Web Service, Connector, SQL...) auf die Information zuzugreifen

---

<sup>3</sup>Application Programming Interface

# Web Scraping

## Die Idee dahinter

- ▶ Informationen aus dem HTML Dokument extrahieren (Achtung:  
Schnell sehr umfangreich!)
- ▶ mit Hilfe von CSS Selektoren die interessanten Information  
identifizieren

# Web Scraping

In R:

- ▶ steht mit **rvest** im tidyverse bereits ein Package bereit

## Review: HTML und CSS Selektoren

# Extracting data from the web

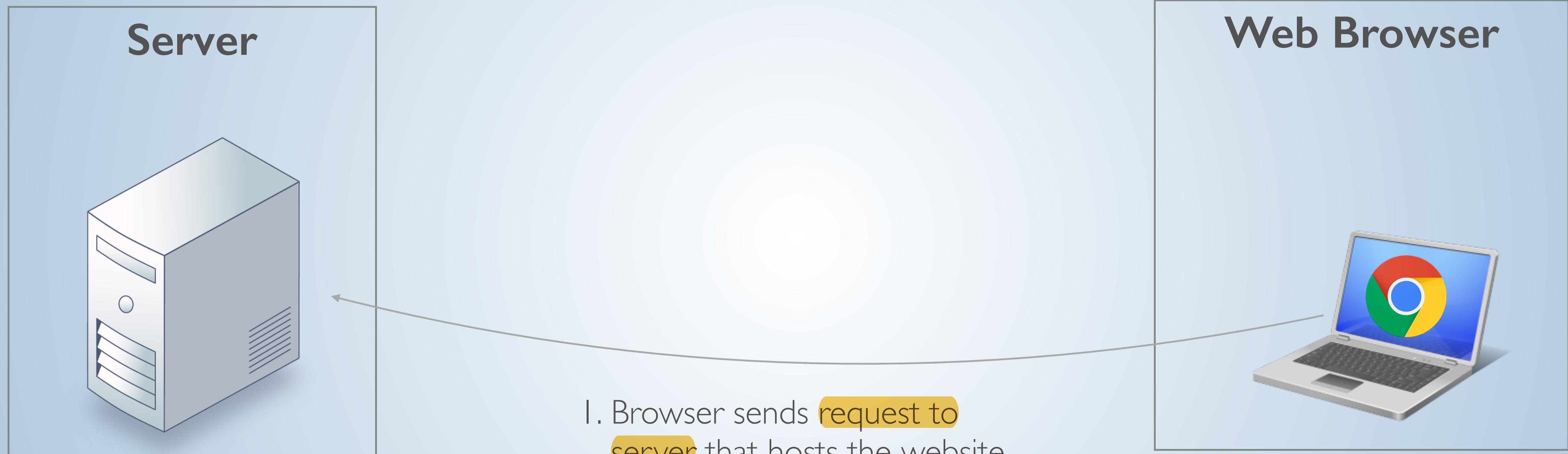
## APIs and beyond



Scott Chamberlain  
Karthik Ram  
**Garrett Grolemund**

June 2016

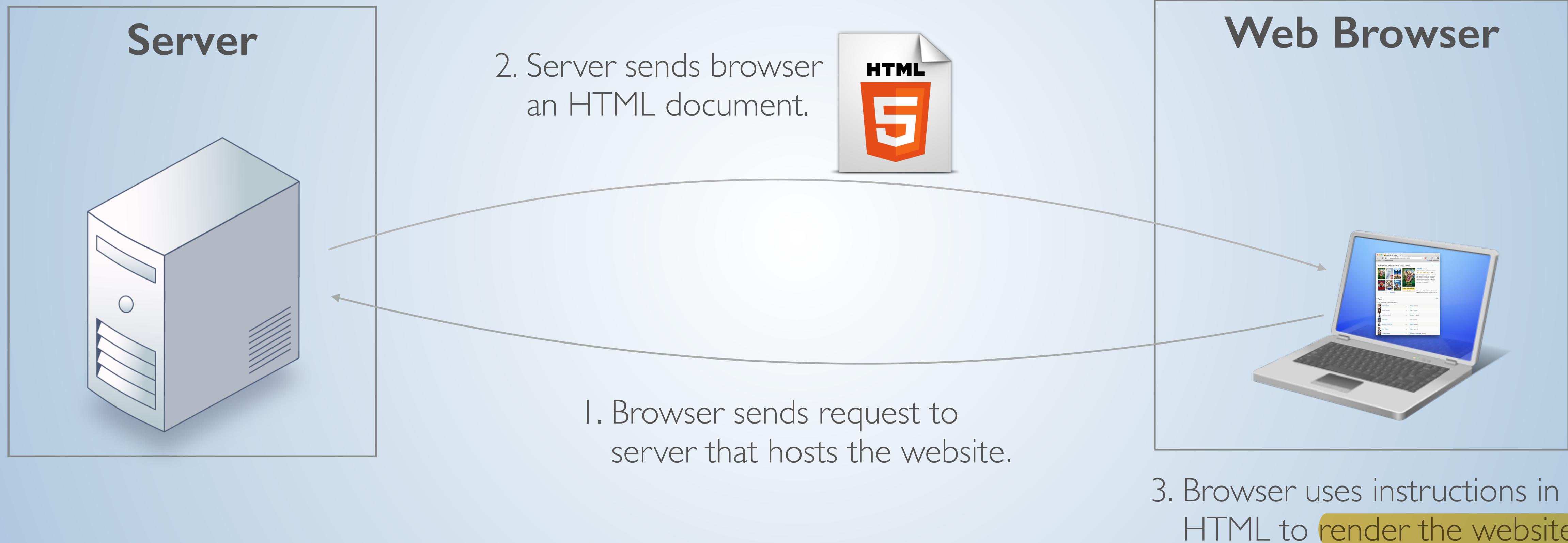
# *HTML (Review)*



# *HTML (Review)*



# *HTML (Review)*

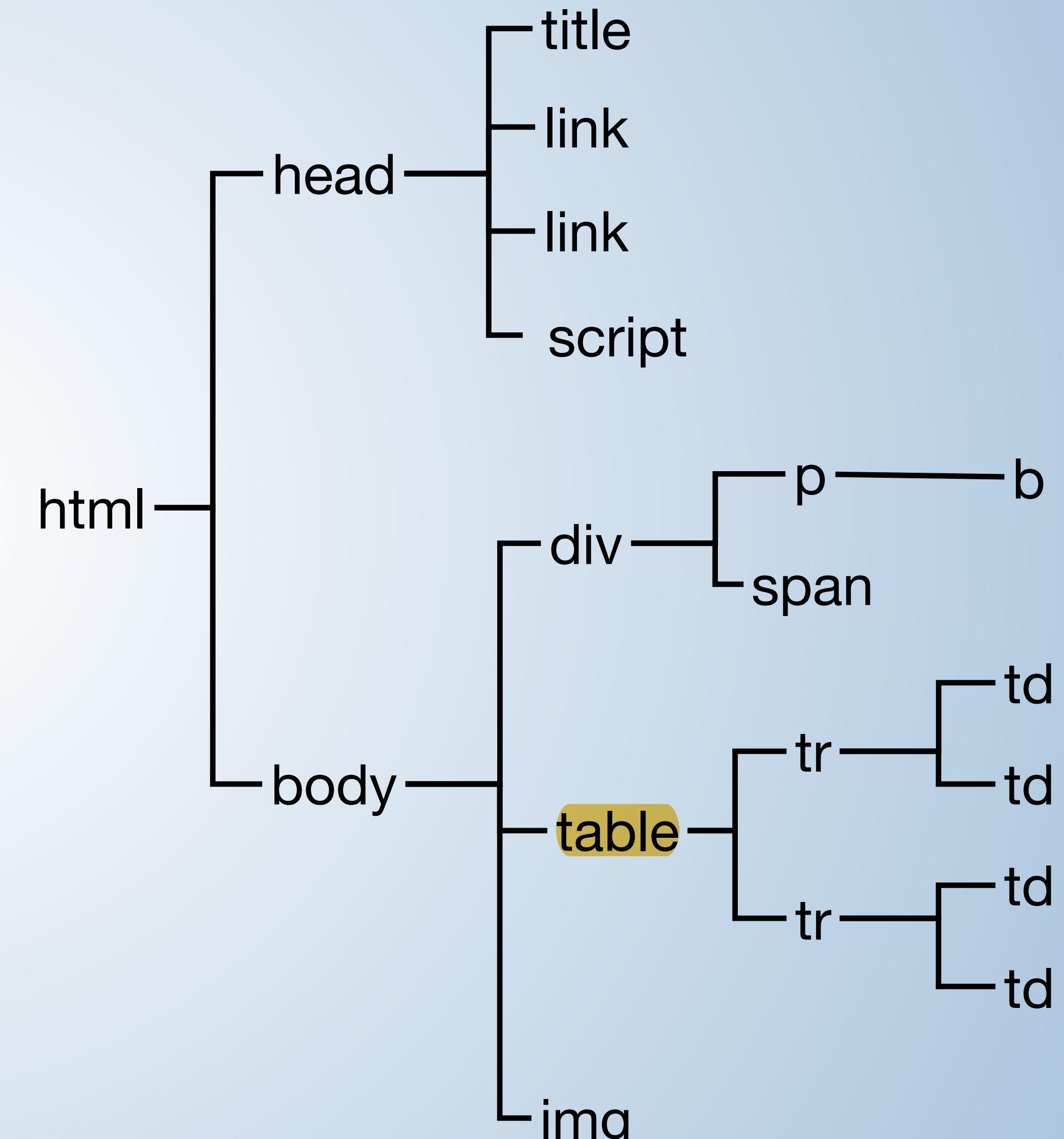


# HTML (Review)



```
<html>
  <head>
    <title>Title</title>
    <link rel="icon" type="icon" href="http://a" />
    <link rel="icon" type="icon" href="http://b" />
    <script type="text/javascript">
      var ue_t0=window.ue_t0||+new Date();
    </script>
  </head>
  <body>
    <div>
      <p>Click <b>here</b> now.</p>
      <span>Frozen</span>
    </div>
    <table style="width:100%">
      <tr>
        <td>Kristen</td>
        <td>Bell</td>
      </tr>
      <tr>
        <td>Idina</td>
        <td>Menzel</td>
      </tr>
    </table>
    
  </body>
</html>
```

HTML Dokument (vereinfacht)



Baumstruktur der HTML Elemente

# HTML (Review)

Each element in the page is created by a **tag**.

```
<a href="http://github.com">GitHub</a>
```

tag name

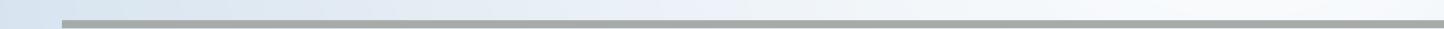
attribute  
(name)

attribute  
(value)

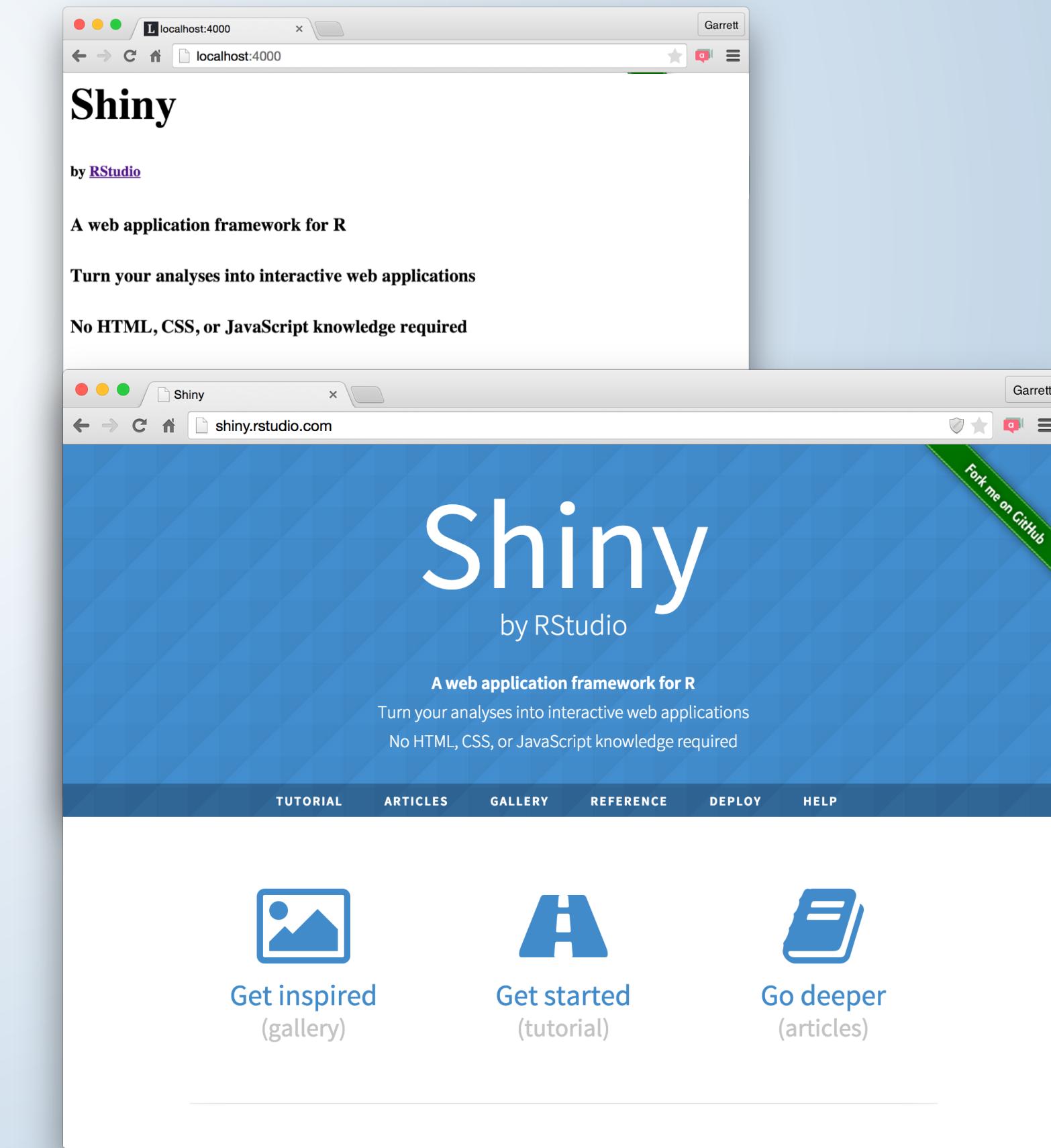
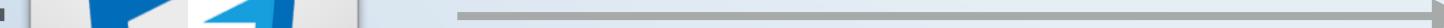
content

# CSS (*Review*)

Cascading Style Sheets (CSS) are a framework for customizing the appearance of elements in a web page.



+



# CSS (*Review*)



localhost:4000 Garrett

## Shiny

by RStudio

A web application framework for R

Turn your analyses into interactive web applications

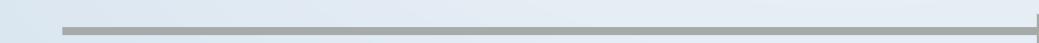
No HTML, CSS, or JavaScript knowledge required

- [Tutorial](#)
- [Articles](#)
- [Gallery](#)
- [Reference](#)
- [Deploy](#)
- [Help](#)

[Get inspired](#)  
(gallery)

[Get started](#)  
(tutorial)

[Go deeper](#)  
(articles)



localhost:4000 Garrett

# Shiny

by RStudio

A web application framework for R

Turn your analyses into interactive web applications

No HTML, CSS, or JavaScript knowledge required

[TUTORIAL](#) [ARTICLES](#) [GALLERY](#) [REFERENCE](#) [DEPLOY](#) [HELP](#)

 [Get inspired](#)  
(gallery)

 [Get started](#)  
(tutorial)

 [Go deeper](#)  
(articles)

Fork me on GitHub

# CSS (Review)



```
span {  
    color: #ffffff;  
}  
  
.num {  
    color: #a8660d;  
}  
  
table.data {  
    width: auto;  
}  
  
#firstname {  
    background-color: yellow;  
}
```

← selector

← styling

← selector

← styling

← selector

← styling

← selector

← styling

Auszug aus einer CSS Datei (vereinfacht)

# CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

tag name

class  
(optional)

id  
(optional)

# CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

span

CSS selector for **ALL** elements with:

- the **span tag**

# CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
.bigname
```

CSS selector for **ALL** elements with:

- the **bigname class**

# CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
span.bigname
```

CSS selector for **ALL** elements with:

- the **span tag**

**AND**

- the **bigname class**

# CSS (*Review*)

A CSS script describes an element by its tag, class, and/or ID.

```
<span class="bigname" id="shiny">Shiny</span>
```

```
#shiny
```

CSS selector for **ALL** elements with:

- the **shiny id**

R Package: rvest

# rvest



A package that makes it easy to extract  
info from a webpage.

```
install.packages("rvest")
```

tidyverse Package

\* This will also install *xml2*, a package that *rvest* relies on.

# Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`

```
library(rvest)  
frozen <- read_html("http://www.imdb.com/title/tt2294629/")
```

read\_html

URL

\* `read_html()` comes in the `xml2` package

# Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`

```
itals <- html_nodes(frozen, "em")
```

XML

hier: emphasized tag

CSS selector

# Basic Workflow

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`
3. Extract content from nodes with `html_text()`,  
`html_name()`, `html_attrs()`, `html_children()`,  
`html_table()`

# Tables

Use `html_table()` to scrape whole tables of data **as  
a data frame.**

```
tables <- html_nodes(kw, css = "table")
html_table(tables, header = TRUE)[[2]]
```

# Recap

1. Download the HTML and turn it into an XML file with `read_html()`
2. Extract specific nodes with `html_nodes()`
3. Extract content from nodes with `html_text()`,  
`html_name()`, `html_attrs()`, `html_children()`,  
`html_table()`

Werkzeug: selectorGadget

# selectorGadget

A GUI tool to identify CSS selector combinations



# To Install

1. Run `vignette("selectorgadget")`
2. Drag `Selectorgadget` link into your browser's bookmark bar

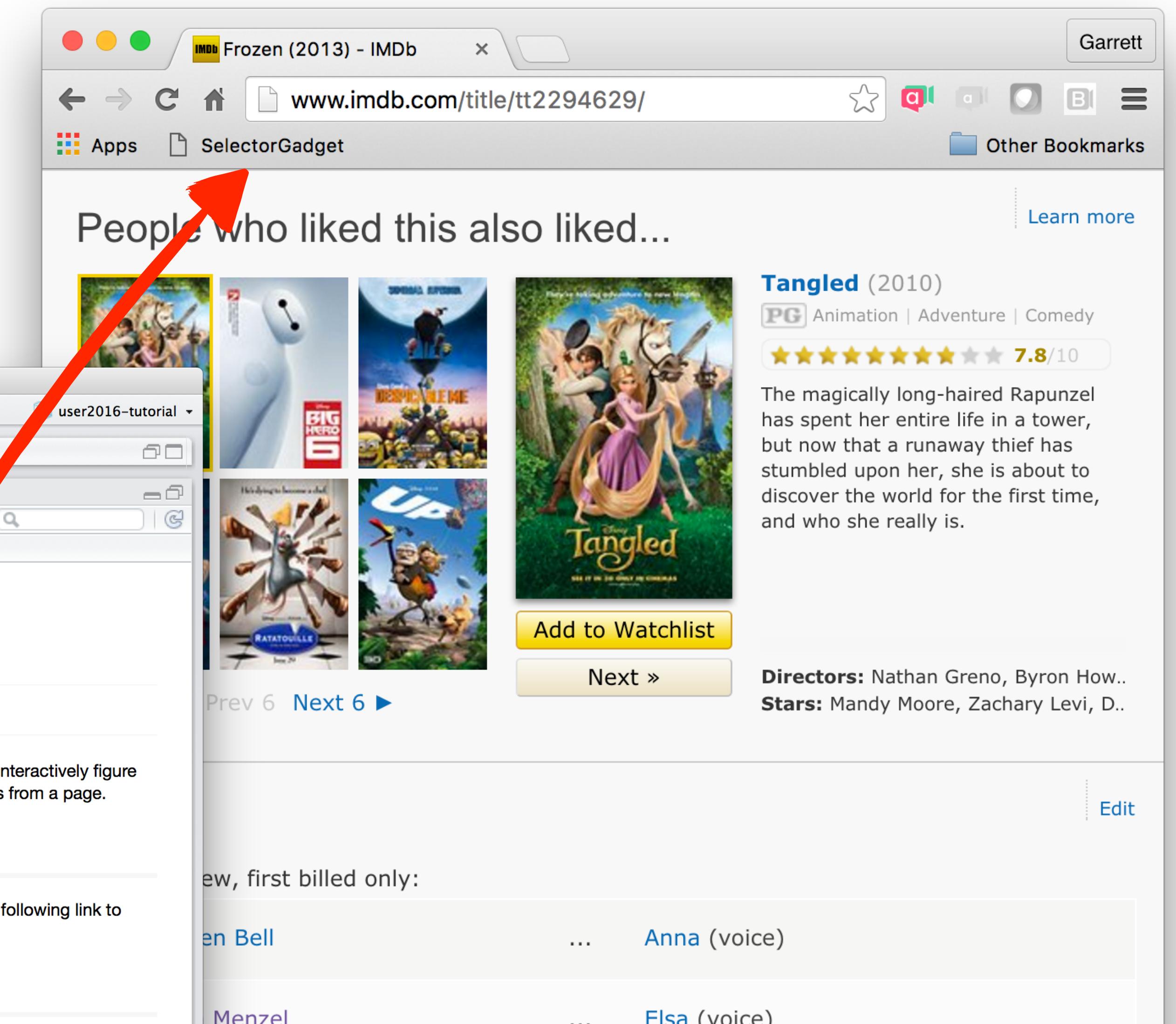
The screenshot shows the RStudio interface. In the top-left pane, there is an R script named "scraping-outline.Rmd" with the following code:

```
1 # Read in web page
2
3 library(rvest)
4 frozen <- read_html("http://www.imdb.com/title/tt2294629")
5
6 # Look at web page
7
8 frozen
9 html_structure(frozen)
10 as_list(frozen)
11
12 xml_children(frozen)
13 xml_children(frozen)[[2]]
```

In the bottom-left pane, the console shows the command:

```
> vignette("selectorgadget")
```

The main pane displays the "Selectorgadget" vignette page. A red circle highlights the instruction: "To install it, open this page in your browser and then drag the following link to your bookmark bar: [Selectorgadget](#)". A red arrow points from this instruction to the "Selectorgadget" link in the browser's address bar.



# To Use

1. **Navigate** to a webpage
2. **Open** the SelectorGadget bookmark
3. **Click** on item to scrape
4. **Click** on yellow items You do not want to scrape
5. **Click** on additional items that you do want to scrape
6. **Copy** selector to use with `html_nodes()`

.fa-bolt

Clear (1)   Toggle Position   XPath   Help   X

CSS selector to use

start over

move gadget

show XPath

help

close gadget

## Rechtliche Bedenken

## Rechtliche Bedenken

- ▶ Urheberrechte
- ▶ Nutzungsbedingungen (Verbot automatischen Auslesens?)
- ▶ ggf. Inhalt der robots.txt Datei prüfen

## Zusammenfassung

## Zusammenfassung

- ▶ Prinzip: Informationen aus HTML extrahieren; dabei HTML/CSS Strukturen nutzen um die interessante Information auf umfangreichen Webseiten zu identifizieren
- ▶ im R Package **rvest** stehen die dazu erforderlichen Funktionen bereit
- ▶ mit dem Bookmarklet **selectorGadget** lassen sich leicht geeignete CSS Selektoren finden
- ▶ immer erst die rechtliche Lage überprüfen

Vielen Dank!