



**OFICINA CIS-USP**

---

**INTRODUÇÃO À PROGRAMAÇÃO E  
WEB SCRAPING**

# O QUE VAMOS APRENDER...

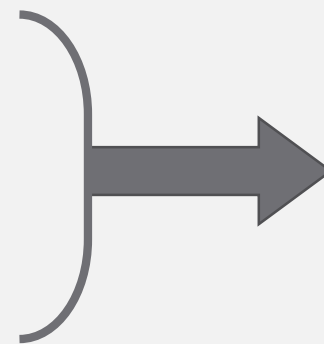
Objetivo: Capacitar pesquisadoras(es) do CIS em coletar dados provenientes da internet por meio de Web Scraping.

- Introdução à programação



**AULA 1**

- Conhecendo alguns pacotes: *magrittr*, *stringr* e *purrr*



**AULA 2**

- Introdução à estrutura HTML

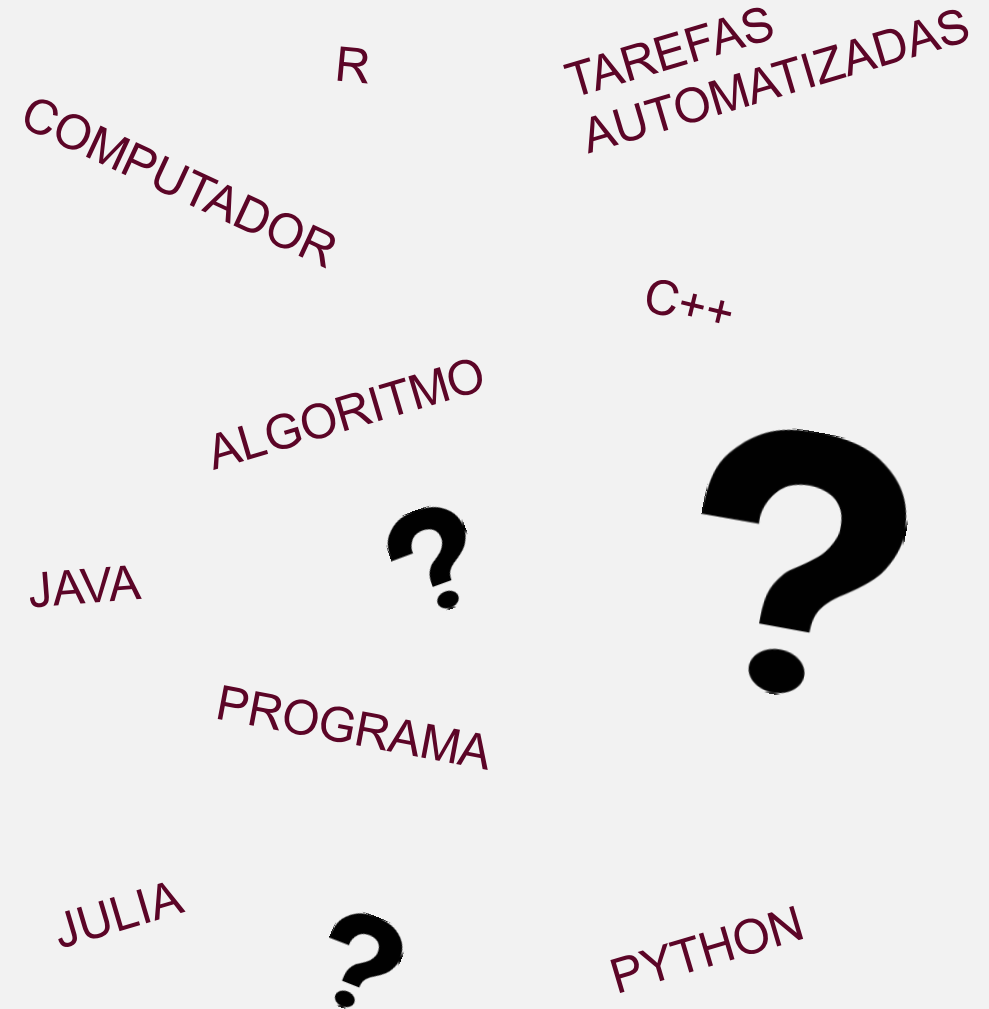
- Web Scraping com *rvest*



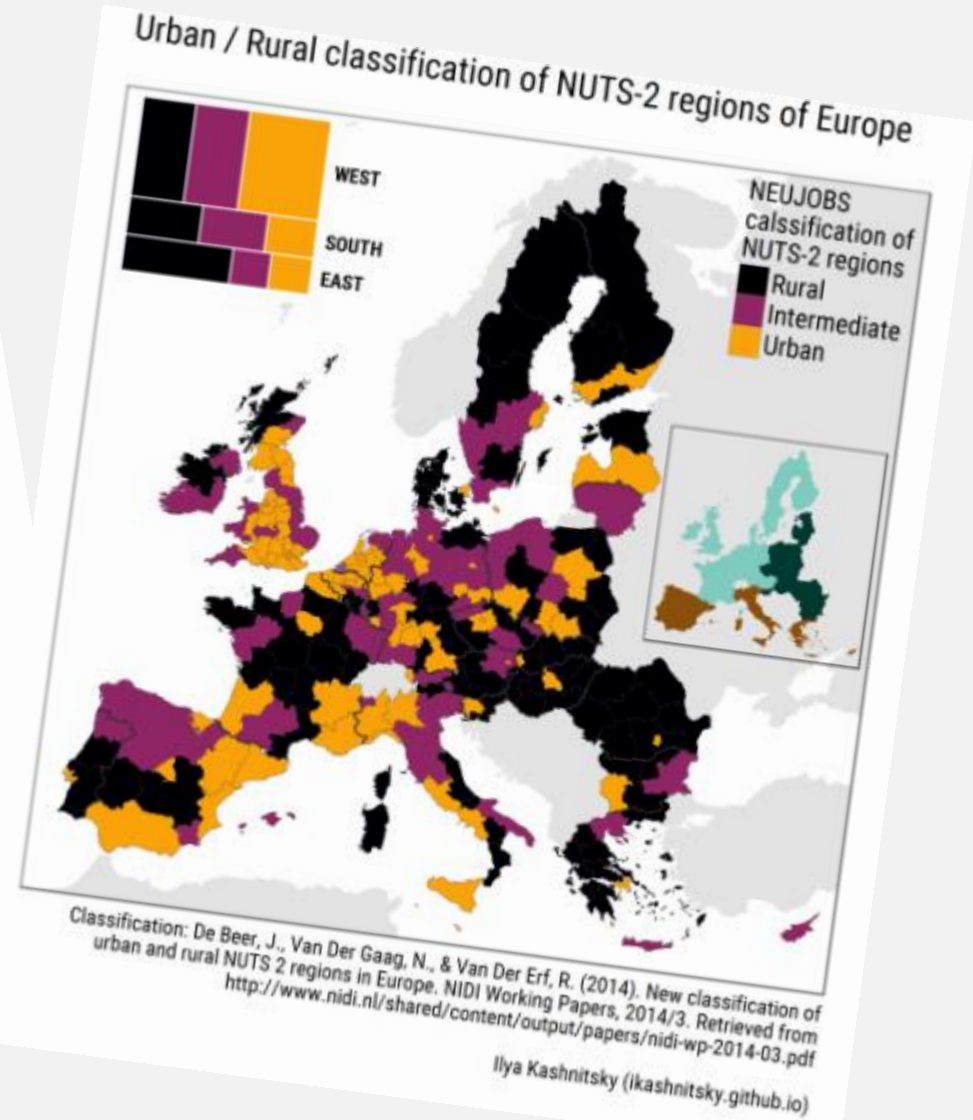
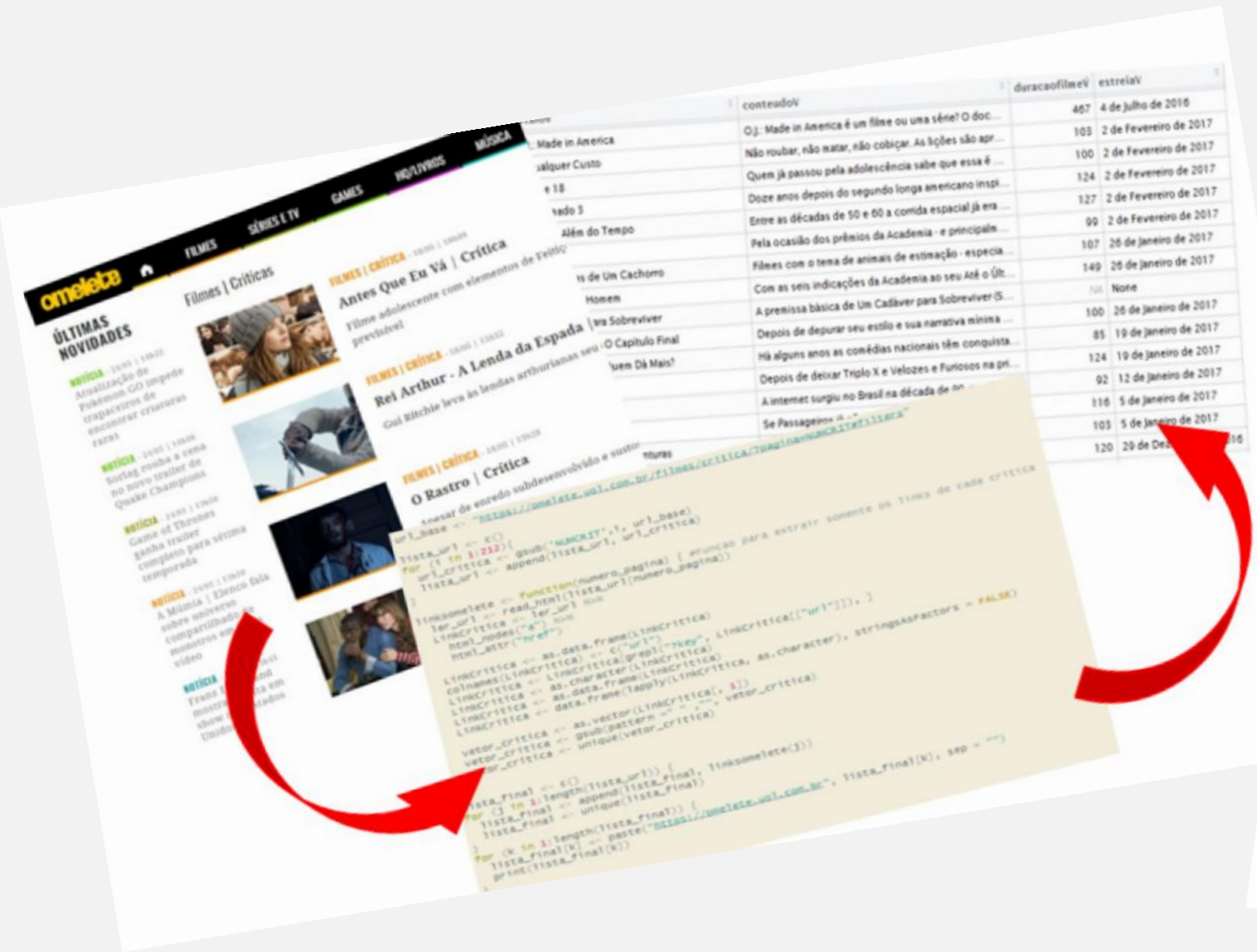
**AULA 3**

# O QUE É PROGRAMAR?

```
11 #Extraindo uma lista com todas as críticas do omelete
12 url_base <- "https://omelete.uol.com.br/filmes/critica/?pagina=NUMCRIT#filters"
13
14 lista_url <- c()
15 for (i in 1:212){
16   url_critica <- gsub('NUMCRIT',i, url_base)
17   lista_url <- append(lista_url, url_critica)
18 }
19
20 linksomelete <- function(numero_pagina) { #Funcao para extrair somente os links de cada critica
21   ler_url <- read_html(lista_url[numero_pagina])
22   LinkCritica <- ler_url %>%
23     html_nodes("a") %>%
24     html_attr("href")
25
26   LinkCritica <- as.data.frame(LinkCritica)
27   colnames(LinkCritica) <- c("url")
28   LinkCritica <- LinkCritica[grepl("?key", LinkCritica[["url"]]), ]
29   LinkCritica <- as.character(LinkCritica)
30   LinkCritica <- as.data.frame(LinkCritica)
31   LinkCritica <- data.frame(lapply(LinkCritica, as.character), stringsAsFactors = FALSE)
32
33   vetor_critica <- as.vector(LinkCritica[, 1])
34   vetor_critica <- gsub(pattern = " ", "", vetor_critica)
35   vetor_critica <- unique(vetor_critica)
36 }
37
38 lista_final <- c()
39 for (j in 1:length(lista_url)) {
40   lista_final <- append(lista_final, linksomelete(j))
```



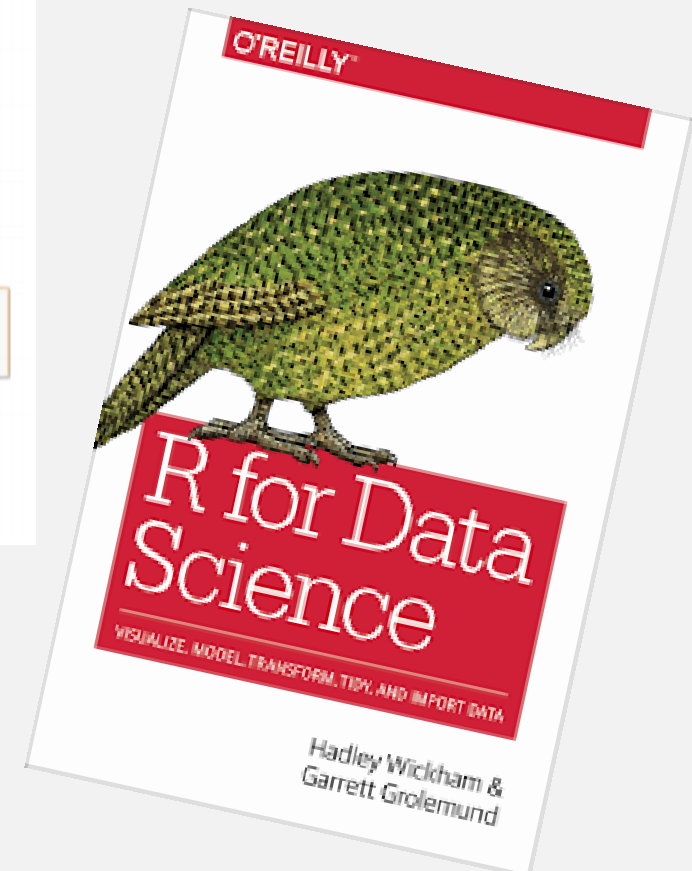
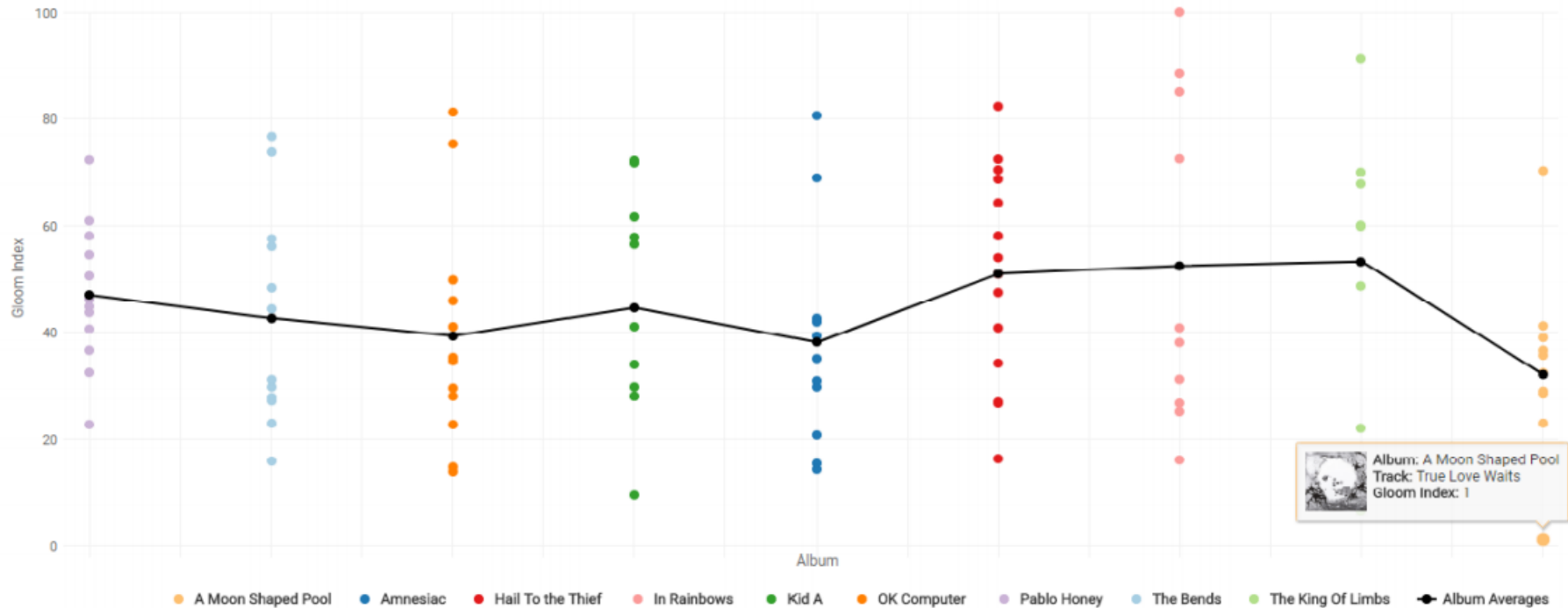
# LINGUAGEM R? O QUE PODEMOS FAZER?



# LINGUAGEM R? O QUE PODEMOS FAZER?

## Data Driven Depression

Radiohead song sadness by album



VAMOS PROGRAMAR!