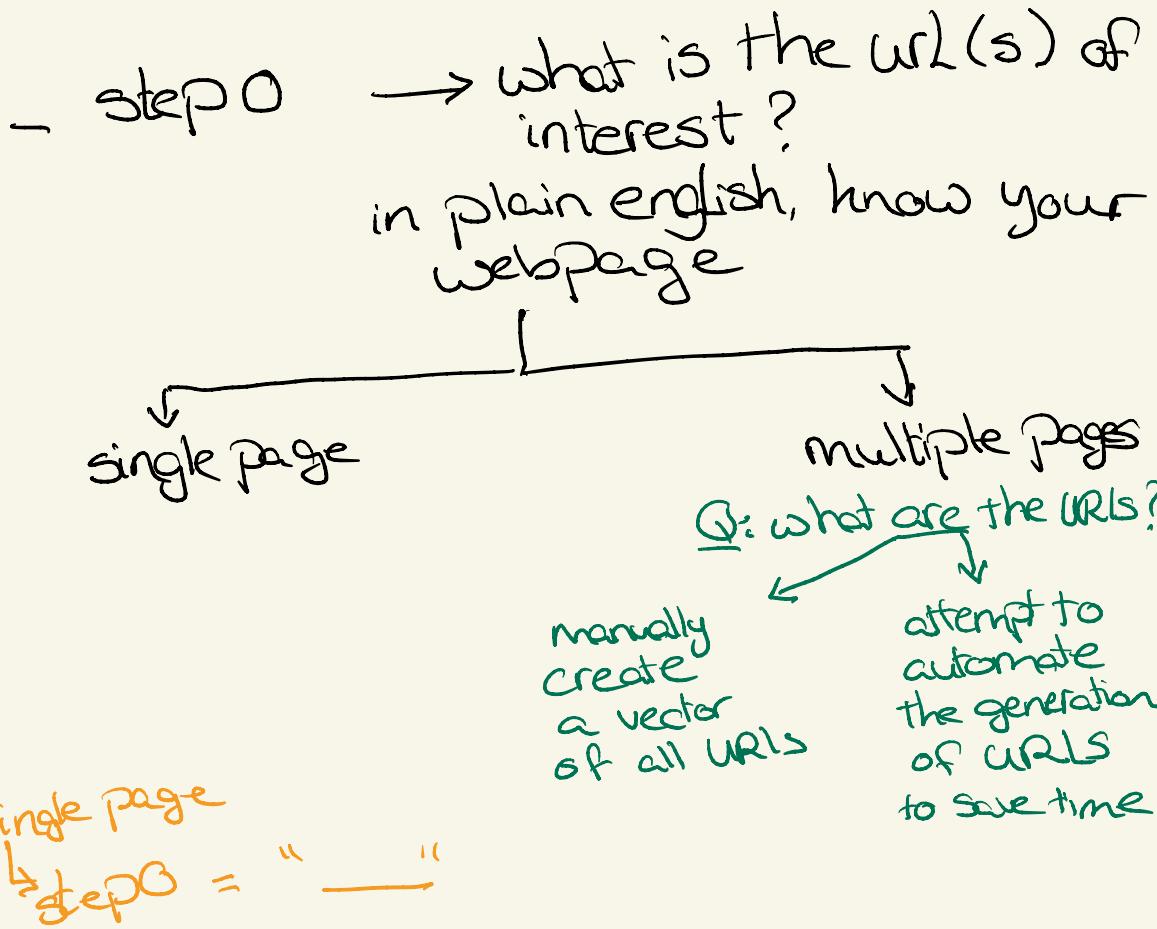


# Web Scraping Process:

## (A) No "Pipe" (%)



multiple pages

increment the URL @ the end (e.g. you will ~~post~~ year or some counter @ the end)

- Step1  $\Rightarrow$  get the html into your software ( backend of the html & Not how your browser is showing it )

Step1 = read\_html( x = url )  
↓  
rvest pkg

For your code to work:



a static webpage



does the website allow you to scrape the data?

if its working

step1 will return a list  
↳ head  
↳ body

- step2 "often repeat that to get diff elements/parts of webpage"

inputs:

- \* read webpage stored in step1
- \* CSS for the elements of int
  - "typically divided by appearance on page"

"iterative"

When to know that your result is correct  
↳ starting, ending & total # results seems to be correct

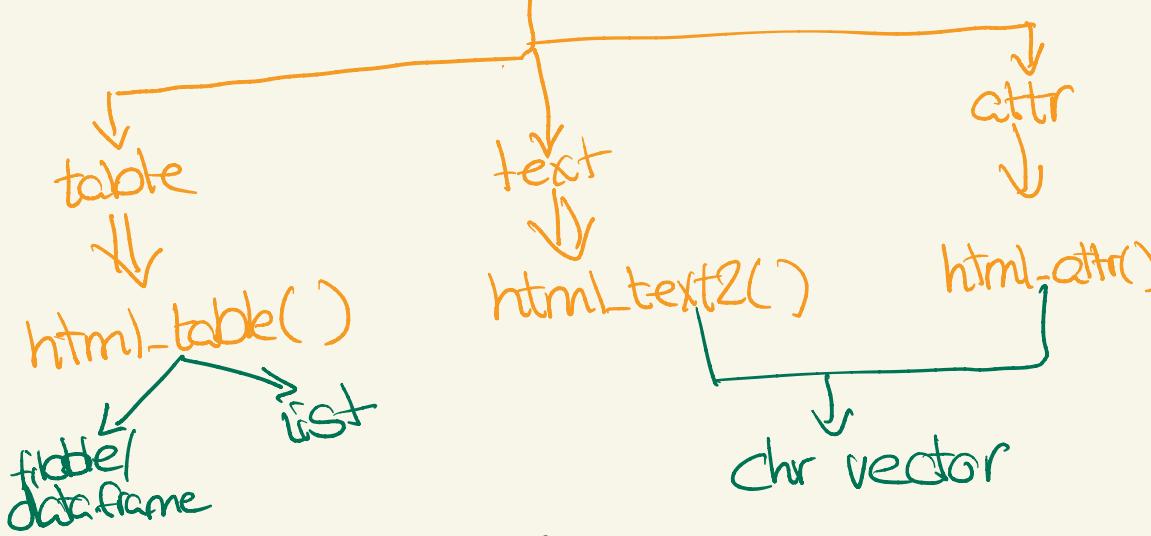
incorrect #'s

→ 1 result

→ see if you have a `nth-child()` in your selector  
→ `html-element()` instead of `html-elements()`

→ too many results ⇒ add tags  
→ too few ⇒ remove classes / tags

- Step 3 → make your results clean  
readable



input: only step 2

output: readable value similar to step 2

- Step 4 → merge diff text and  
(not always) put in a tibble  
needed)

to make it work: Length of the  
diff objs has to be the same

if you pipe, you do not need  
to store intermediate results.

(you do not have to pipe  
if you don't like it)

native pipe: |>

same as the magrittr/dplyr  
pipe %>%