# CIÊNCIA DE DADOS (BIG DATA PROCESSING AND ANALYTICS)

RECUPERAÇÃO DA INFORMAÇÃO NA WEB E EM REDES SOCIAIS





# TRILHA 2 RECUPERAÇÃO DE INFORMAÇÃO POR RASPAGEM - INTRODUÇÃO

# Qual é a diferença entre internet e web?

#### **Internet**

conexão • rede de rede • tcp/ip

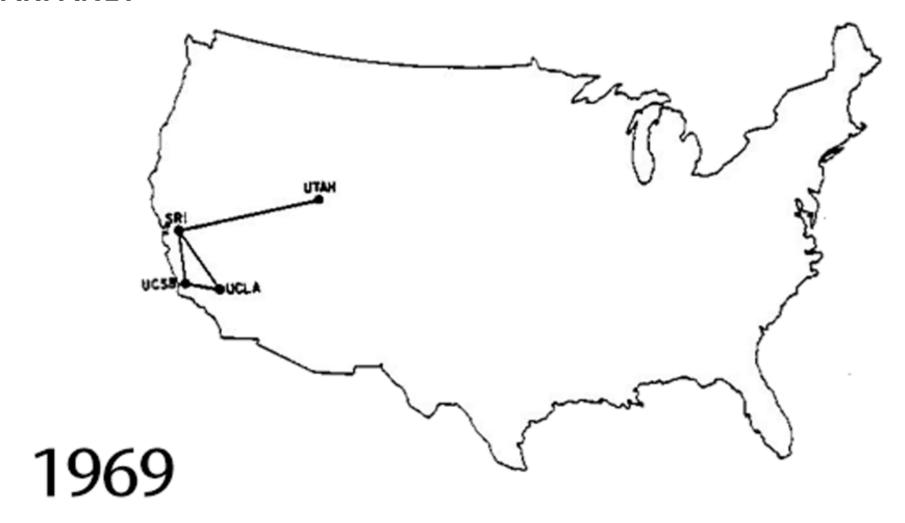


#### **World Wide Web**

aplicação • distribuição de conteúdo • http



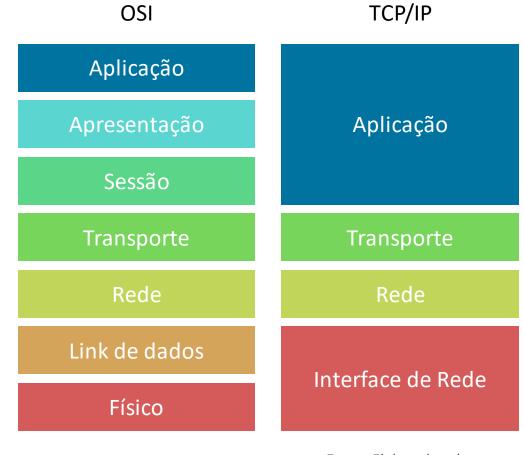
## **ARPANET**



ROSER, Max; RITCHIE, Hannah; ORTIZ-OSPINA, Esteban. Internet. Our WorldInData.org, 2015. Disponível em: <a href="https://ourworldindata.org/internet">https://ourworldindata.org/internet</a>.

# Algumas aplicações famosas da internet

- File Transfer Protocol (FTP) e SSH File Transfer Protocol (SFTP)
- Simple Mail Transfer Protocol (SMTP), POP e IMAP
- Network News Transfer Protocol (NNTP)
- Internet Relay Chat (IRC)
- Gopher
- HTTP



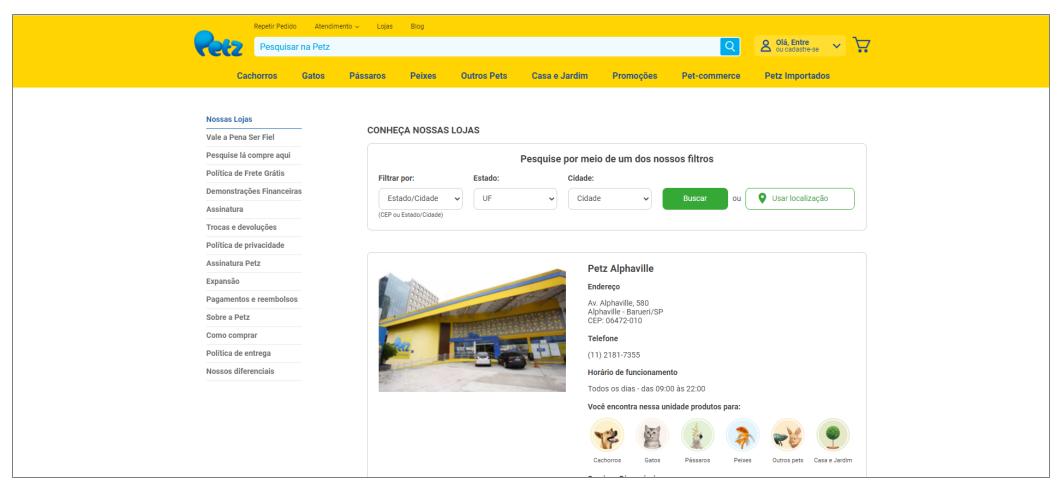
Fonte: Elaborada pelo autor.

## A web como nós a conhecemos

### **HTTP**

**HTML** • CSS • Javascript • Imagens • Vídeos • etc.

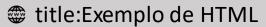
# Uma página da internet



**Fonte:** <a href="https://www.petz.com.br/nossas-lojas">https://www.petz.com.br/nossas-lojas</a>.

## Uma página da internet – Olá, mundo!

```
♦ hello.html > ♦ html
     <!DOCTYPE html>
     <html>
         <head>
            <title>Exemplo de html</title>
         </head>
         <body>
            <h1>Exemplo de HTML</h1>
 7
            <h2>0 classico Olá mundo</h2>
 8
            Este paragrafo vai aparecer abaixo do <i>heading</i> nivel 2
 9
            <h2>Exemplo com lista</h2>
10
            (ul>
11
                Item 1 
12
                Item 2 
13
                Item 3 
14
15
            </body>
16
     </html>
17
```



### Exemplo de HTML

#### O classico Olá mundo

Este paragrafo vai aparecer abaixo do heading nivel 2 de Olá Mundo

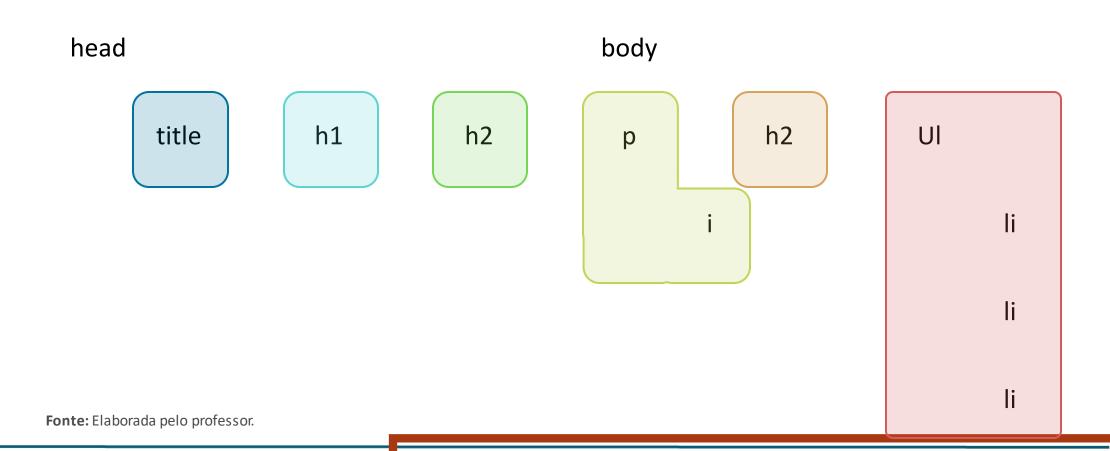
#### Exemplo com lista

- Item 1
- Item 2
- Item 3

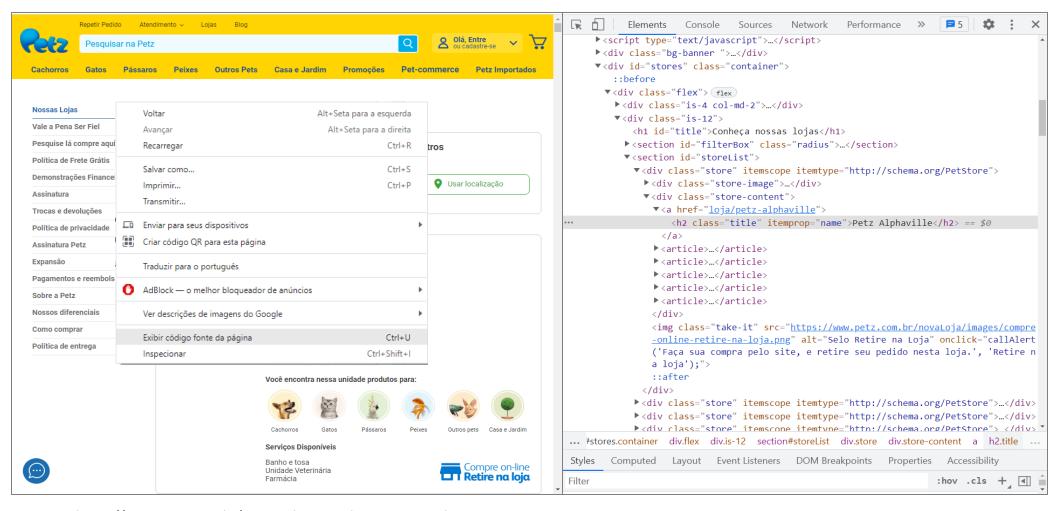
**Fonte:** Elaborada pelo professor.

# Uma página da internet – Olá, mundo! – Árvore DOM

document



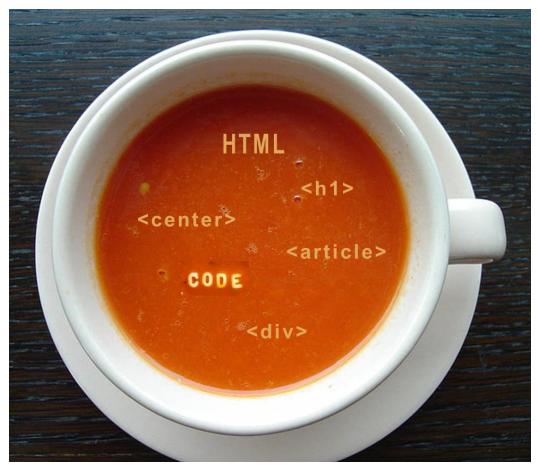
## Uma página da internet e seu HTML



Fonte: <a href="https://www.petz.com.br/nossas-lojas">https://www.petz.com.br/nossas-lojas</a>; Chrome DevTools.

# A biblioteca Beautiful Soup 4.0

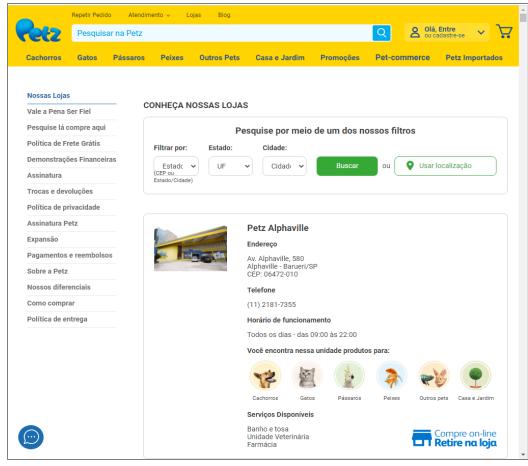
- Muitos documentos possuem erro de sintaxe (sopa de etiquetas).
- Leonard Richardson criou um interpretador que trata alguns dos erros mais comuns.
- Ideal para pequenos projetos de raspagem de documentos.



Adaptação em foto de domínio público. **Fonte:** <a href="https://www.hippopx.com/">https://www.hippopx.com/>.</a>

## Exercício prático

equipe de marketing de uma empresa de distribuição de produtos veterinários precisava do endereço de todas as lojas da rede Petz para uma campanha. Sabe-se que todos os endereços estão disponíveis no porém, como empresa, utilizar a raspagem de tela e Beautiful Soup 4.0 para obter esses dados de forma estruturada?



**Fonte:** <a href="https://www.petz.com.br/nossas-lojas">https://www.petz.com.br/nossas-lojas</a>.

# Passo 1 – Baixando o documento para uma *string python*

```
    import requests
    url = 'https://www.petz.com.br/nossas-lojas'
    res = requests.get(url)
    html_page = res.text
```

# Passo 2 – Invocando o *Beautiful Soup* 4.0

```
1. import requests
2.
3. url = 'https://www.petz.com.br/nossas-lojas'
4. res = requests.get(url)
5. html_page = res.text

1. from bs4 import BeautifulSoup
2. soup = BeautifulSoup(html_page, 'html.parser')
```

## Passo 3 – Fazendo a raspagem

address = address.get text()

row = { 'name': name,

stores.append(row)

address = re.sub(' +', ' ', address)

'telephone': telephone.

'addressLocality': addressLoc,

6.

7.

8.

9.

10.

11.

12.

13.

14. 15.

16.

address = " ".join([s for s in address.strip().splitlines(False) if s.strip()])

region = store.find("span", {"itemprop": "addressRegion"}).get text().strip()

openingHours = store.find("p", {"itemprop": "openingHours"}).get text().strip()

telephone = store.find("p", {"itemprop": "telephone"}).get\_text().strip()

addressLoc = store.find("span", {"itemprop": "addressLocality"}).get\_text().strip()

'address':address,

'region': region}

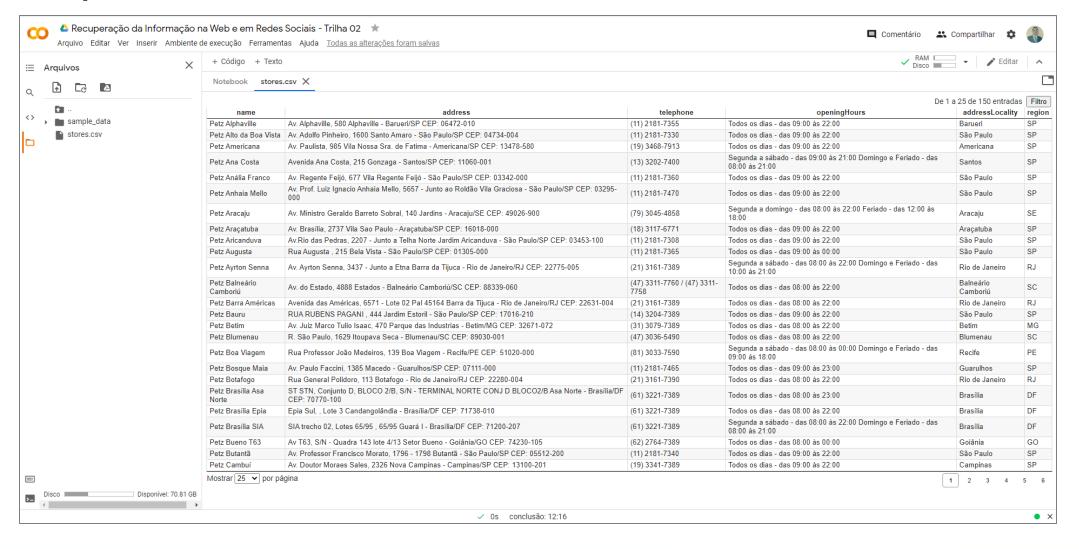
'openingHours':openingHours,

## Passo 4 – Exportando para o Excel

dict writer.writerows(stores)

```
1. import requests
3. url = 'https://www.petz.com.br/nossas-lojas'
4. res = requests.get(url)
html page = res.text
1. from bs4 import BeautifulSoup
2. soup = BeautifulSoup(html page, 'html.parser')
1. import re
2. stores = []
3. for index, store in enumerate(soup.find_all("div", {"class": "store"}), start=1):
       name = store.find("h2", {"itemprop": "name"}).get text().strip()
4.
       address = store.find("p", {"itemprop": "address"})
5.
       address = address.get text()
6.
       address = re.sub(' +', ' ', address)
7.
8.
       address = " ".join([s for s in address.strip().splitlines(False) if s.strip()])
       addressLoc = store.find("span", {"itemprop": "addressLocality"}).get_text().strip()
9.
       region = store.find("span", {"itemprop": "addressRegion"}).get_text().strip()
10.
       telephone = store.find("p", {"itemprop": "telephone"}).get_text().strip()
11.
       openingHours = store.find("p", {"itemprop": "openingHours"}).get_text().strip()
12.
13.
       row = { 'name': name,
                                                     'address':address,
               'telephone': telephone.
                                                     'openingHours':openingHours,
14.
               'addressLocality': addressLoc,
                                                     'region': region}
15.
16.
       stores.append(row)
1. import csv
2. keys = stores[0].keys()
3. with open('stores.csv', 'w', newline='') as output file:
      dict writer = csv.DictWriter(output_file, keys)
4.
     dict writer.writeheader()
5.
```

## Arquivo final tabulado



# Disponível no Google colab



https://bit.ly/3yOaPJM



