

PokeCrawler

Ahmed Mohamad Bakri
Matheus Miranda

PokeCrawler

Neste trabalho implementamos um crawler que rasteja através de uma url que no caso é a de uma pokédex para extrairmos dados como id, nome, habilidade, evoluções e tipos



Como

Através de seletores css que pegamos inspecionando elemento nas partes do site que desejamos como o id de um pokémon, e utilizando o método response do scrapy nós conseguimos pegar o conteúdo desses seletores



Complete Pokémon Pokédex

This is a full list of every Pokémon from all 9 generations of the Pokémon series, along with their main stats.

The table is sortable by clicking a column header, and searchable by using the controls above it.

Name: Type: - All -

#	Name	Type	Total	HP	Attack	Defense	Sp. Atk	Sp. Def	Speed
0001	Bulbasaur	GRASS POISON	318	45	49	49	65	65	45
0002	Ivysaur	GRASS POISON	405	60	62	63	80	80	60
0003	Venusaur	GRASS POISON	525	80	82	83	100	100	80
0003	Venusaur Mega Venusaur	GRASS POISON	625	80	100	123	122	120	80
0004	Charmander	FIRE	309	39	52	43	60	50	65
0005	Charmeleon	FIRE	405	58	64	58	80	65	80

DevTools is now available in Portuguese!

Always match Chrome's language Switch DevTools to Portuguese Don't show again

Elements Console Sources Network

```
<div class="resp-scroll">
  <table id="pokedex" class="data-table sticky-header block-wide" style="opacity: 1;">
    <thead>
      <tr>
        <td class="cell-num cell-fixed" data-sort-value="1">
          <td class="cell-name">
            <a class="ent-name" href="/pokedex/bulbasaur" title="View Pok
              edex for #0001 Bulbasaur">Bulbasaur</a> == $0
            </td>
          <td class="cell-icon">
            <td class="cell-num cell-total">318</td>
            <td class="cell-num">45</td>
            <td class="cell-num">49</td>
            <td class="cell-num">49</td>
            <td class="cell-num">65</td>
            <td class="cell-num">65</td>
          </tr>
        </thead>
        <tbody>
          <tr>
            <td class="cell-num">405</td>
            <td class="cell-num">60</td>
            <td class="cell-num">62</td>
            <td class="cell-num">63</td>
            <td class="cell-num">80</td>
            <td class="cell-num">80</td>
            <td class="cell-num">60</td>
          </tr>
          <tr>
            <td class="cell-num">525</td>
            <td class="cell-num">80</td>
            <td class="cell-num">82</td>
            <td class="cell-num">83</td>
            <td class="cell-num">100</td>
            <td class="cell-num">100</td>
            <td class="cell-num">80</td>
          </tr>
          <tr>
            <td class="cell-num">625</td>
            <td class="cell-num">80</td>
            <td class="cell-num">100</td>
            <td class="cell-num">123</td>
            <td class="cell-num">122</td>
            <td class="cell-num">120</td>
            <td class="cell-num">80</td>
          </tr>
          <tr>
            <td class="cell-num">309</td>
            <td class="cell-num">39</td>
            <td class="cell-num">52</td>
            <td class="cell-num">43</td>
            <td class="cell-num">60</td>
            <td class="cell-num">50</td>
            <td class="cell-num">65</td>
          </tr>
          <tr>
            <td class="cell-num">405</td>
            <td class="cell-num">58</td>
            <td class="cell-num">64</td>
            <td class="cell-num">58</td>
            <td class="cell-num">80</td>
            <td class="cell-num">65</td>
            <td class="cell-num">80</td>
          </tr>
        </tbody>
      </table>
    </div>
```

table#pokedex.data-table.sticky-header.block-wide tbody tr td.cell-name a.ent-name

Styles Computed Layout Event Listeners DOM Breakpoints Properties

Filter .hov .cls + -

```
element.style {
}

.ent-name {
  font-size: 1rem;
  font-weight: bold;
}

a {
  color: #279b9e;
  text-decoration: none;
  transition: color .2s, background-color .2s;
}

a {
  -webkit-box-sizing: border-box;
  box-sizing: border-box;
}

a:webkit-any-link {
}
```

Console What's New

Highlights from the Chrome 116 update

Improved debugging of missing stylesheets

Find and fix issues with missing stylesheets with ease.

Add attribute

Edit attribute

Edit as HTML

Duplicate element

Delete element

Cut

Copy

Paste

Hide element

Force state

Break on

Expand recursively

Collapse children

Capture node screenshot

Scroll into view

Focus

Badge settings...

Store as global variable

Copy element

Copy outerHTML

Copy selector

Copy JS path

Copy styles

Copy XPath

Copy full XPath

Como

No último slide há um exemplo de como o seletor css do nome de um pokémon é retirado.



Exemplo

Este é o nosso seletor do nome do pokémon

```
name_s = "#main > h1::text"
```

E aqui nós pegamos o valor atrelado a esse seletor

```
name = response.css(name_s)
```

através de um yield nós “jogamos” esses dados em csv ou json

```
yield {'Id': id.get(), 'Nome': name.get(), 'Type': type.getall(), 'Altura': height.get(), 'Peso': weight.get(), "Urls das Habilidades": ability_url, "Nomes di
```

Exemplo

Esse é um exemplo do resultado do crawler para um único pokémon no caso o Bulbasaur

```
{"Id": "0001", "Nome": "Bulbasaur", "Type": ["Grass", "Poison", "Cerulean City"],  
"Altura": "0.7\u00a0m (2\u2032\u2033)", "Peso": "6.9\u00a0kg (15.2\u00a0lbs)",  
"Urls das Habilidades": ["/ability/overgrow", "/ability/chlorophyll"], "Nomes das  
Habilidades": ["overgrow", "chlorophyll"], "Evolu\u00e7\u00f5es": ["Bulbasaur",  
"Ivysaur", "Venusaur"]},
```



Tratamento de dados

O nome das habilidades pegamos através das urls das habilidades de um dado pokémon, para isso splitamos a / para pegar as palavras e adicionamos em um array só as palavras relacionados ao nome delas

```
Hab = []
for i in range(len(ability_url)):
    Hab.append(ability_url[i].split("/"))

nameHab = []
for i in range(len(Hab)):
    nameHab.append(Hab[i][2])
```