



UNIVERSIDADE DO MINHO
MESTRADO EM ENGENHARIA INFORMÁTICA

SCRIPTING NO PROCESSAMENTO DE LINGUAGEM
NATURAL

Beautiful Soup - Python

Grupo:

João Vieira	A76516
Manuel Monteiro	A74036
Miguel Dias	PG41089

4 de Maio de 2020

Conteúdo

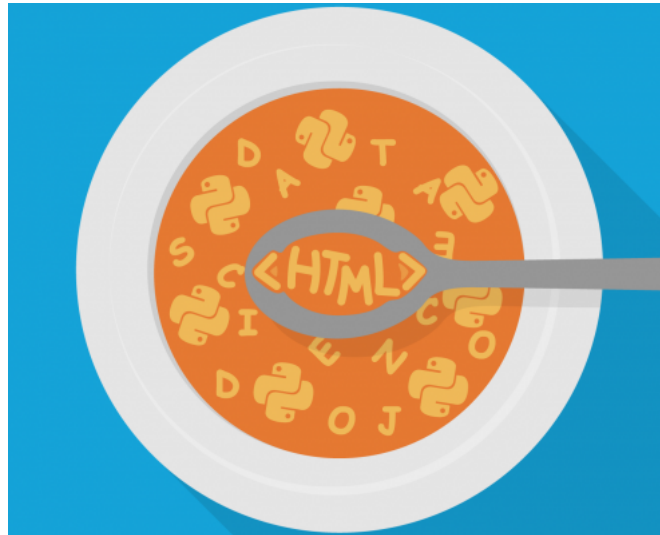
1	Introdução	2
2	Objetivo do Projeto	2
3	Resolução	2
4	Conclusão	8

1 Introdução

Este relatório é relativo a um trabalho realizado para a unidade curricular de *Scripting* no Processamento de Linguagem Natural. Foi escolhida uma ferramenta de uma lista dada pelo professor, sendo escolhida por nós o *Beautiful Soup*. Esta ferramenta de *Python* serve para extrair dados de ficheiros *HTML/XML*. Com ele podemos navegar e pesquisar e modificar a árvore de *parsing* do documento.

2 Objetivo do Projeto

O objetivo deste projeto é utilizar a ferramenta (*Beautiful Soup*) para criar um *script*. Como é sabido, temos outro projeto relativo à unidade curricular do perfil de mestrado. Como tal, decidimos utilizar este projeto para ser utilizado nesse projeto, portanto foi elaborado um *script* para ser utilizado na povoação/atualização dos dados (preço, promoções...) de jogos da plataforma *Steam*.



3 Resolução

Para tal, fazemos um pedido a uma *webpage* que tem uma tabela com todas as promoções existentes na plataforma (Steam DB), a partir do qual utilizamos a ferramenta para obtermos o preço e as promoções dos mesmos.

```

1 <tr class="app appimg" data-appid="1600" data-cache="1587490676">
2   <td data-sort="30">
3     <a target="_blank" rel="noopener"
4       ↪ href="https://store.steampowered.com/app/1600/"
5       ↪ class="info-icon"
6       ↪ title="Store"></a>
7   </td>
8   <td class="applogo">
9     <a target="_blank" rel="noopener" href="/app/1600/"
10    ↪ tabindex="-1" aria-hidden="true">
11      
12    </a>
13  </td>
14  <td>
15    <a target="_blank" rel="noopener" class="b"
16    ↪ href="/app/1600/">Dangerous Waters</a>
17    <span class="subinfo">
18      <span class="highest-discount">lowest price is
19      ↪ <b>1,49</b> at -90%</span>
20    </span>
21  </td>
22  <td class="" data-sort="80">-80%</td>
23  <td data-sort="559">5,59</td>
24  <td data-sort="70.43">70.43%</td>
25  <td data-sort="1588611600" class="timeago"></td>
26  <td data-sort="1587403211" class="timeago"></td>
27  <td data-sort="1139270400" class="timeago"></td>
28 </tr>

```

1: *Excerto HTML da tabela*

```

1 http = urllib3.PoolManager()
2 r = http.request("GET", "https://steamdb.info/sales/")
3
4 soup = BeautifulSoup(r.data, "lxml")
5
6 genRe = r'Genre: [\s]+([\w ,]+) '
7 devRe = r'Developer: [\s]+([\w ,]+) '
8 pubRe = r'Publisher: [\s]+([\w ,]+) '
9 rdRe = r'Release Date: [\s]+([\w ,\d]+) '
10

```

```

11 f = open('dataset.json','w');
12 f.write("[")
13 counter = 0
14
15 page = soup.tbody.find_all("tr")
16 for row in page:
17
18     counter = counter+1
19     print(counter)
20
21     td_list = row.find_all("td")
22     game_page = td_list[0].a["href"]
23
24     name = td_list[2].a.text
25
26     discount = td_list[3].text
27     price = td_list[4].text
28     rating = td_list[5].text
29
30     time_now = dt.datetime.now()
31     sale_end =
        ↳ dt.datetime.fromtimestamp(float(td_list[6]["data-sort"]))
        ↳ - time_now
32     sale_start = time_now -
        ↳ dt.datetime.fromtimestamp(float(td_list[7]["data-sort"]))

```

2: Script que extrai informação da tabela

Além disso para cada jogo encontrado fazemos um outro pedido para obter mais informações sobre o jogo tais como a sua descrição, género, data de lançamento entre outros. Este último pedido é efetuado na página da loja de cada jogo baseado na (Loja da Steam).

```

<meta name="Description" content="Counter-Strike: Global Offensive
↳ (CS: GO)
expands upon the team-based action gameplay that it pioneered when
↳ it was launched 19 years ago.
CS: GO features new maps, characters, weapons, and game modes,
and delivers updated versions of the classic CS content (de_dust2,
↳ etc.).">

<div class="details_block">
  <b>Título:</b> Counter-Strike: Global Offensive<br>
  <b>Gênero:</b>
  <a href="...">Ação</a>,
  <a href="...">
    Grátis para Jogar</a>
  <br>
  <div class="dev_row">
    <b>Desenvolvedor:</b>
    <a href="...">
      Valve
    </a>,
    <a href="...">
      Hidden Path Entertainment</a>
  </div>
  <div class="dev_row">
    <b>Editora:</b>
    <a href="...">Valve</a>
  </div>
  <b>Data de lançamento:</b> 21 Ago, 2012<br>
</div>

```

3: Informações do jogo

```

1 # PÁGINA DO JOGO
2 g = requests.get(game_page)
3
4 if(g.status_code == 200):
5
6     game = BeautifulSoup(g.content, "lxml")
7     gdesc = game.find("meta",{'name':
↳ "Description"})["content"]

```

```

8      gtable = game.find("div",{ 'class': "details_block"})
9
10     if(gtable):
11         genre = re.findall(genRe,gtable.text)
12         if(genre):
13             genrel = map(lambda x :
14                 ↪ x.strip(),genre[0].split(","))
15         else:
16             genrel = []
17
18         developer = re.findall(devRe,gtable.text)
19         if(developer):
20             developerl = map(lambda x :
21                 ↪ x.strip(),developer[0].split(","))
22         else:
23             developerl = []
24
25         publisher = re.findall(pubRe,gtable.text)
26         if(publisher):
27             publisherl = map(lambda x :
28                 ↪ x.strip(),publisher[0].split(","))
29         else:
30             publisherl = []
31
32         releaseDate = re.findall(rdRe,gtable.text)
33         if (releaseDate):
34             releaseDate = releaseDate[0]
35         else:
36             releaseDate = ""

```

4: Script que extrai a informação de um jogo

Por fim, é criado um ficheiro em formato *JSON* que guarda toda a informação extraída.

```
1 #File write
2     f.write("{\n")
3     f.write(f'"name": "{name}",\n') #Nome
4     f.write(f'"description": "{gdesc}",\n') #Descrição
5     f.write(f'"discount": "{discount}",\n') #Desconto
6     f.write(f'"price": "{price}",\n') #Preço
7     f.write(f'"rating": "{rating}",\n') #Avaliação
8     f.write(f'"sale_start": "{sale_start}",\n') #Início de
9     ↪ promoção
10    f.write(f'"sale_end": "{sale_end}",\n') #Fim de
11    ↪ promoção
12    f.write(f'"genres": {json.dumps(list(genrel))},\n')
13    ↪ #Generos
14    f.write(f'"developers":
15    ↪ {json.dumps(list(developerl))},\n') #Developers
16    f.write(f'"publishers":
17    ↪ {json.dumps(list(publisherl))},\n') #Publisher
18    f.write(f'"release_date": "{releaseDate}"\n') #Data de
19    ↪ Lançamento
20    if (counter < len(page)):
21        f.write("},\n")
22    else:
23        f.write("}\n")
24
25 f.write("]")
26 f.close()
```

5: Script que guarda a informação no ficheiro

4 Conclusão

A resolução deste pequeno projeto e a utilização do *Beautiful Soup*, permitiu com que nos fosse introduzido o conceito de *Web Scraping*. O módulo estudado fez com que o conceito fosse bem consolidado, devido à facilidade que este dá ao programador de extrair informações de qualquer sítio da *web*. A ferramenta permite criar *script's* que extraem muita informação e ao mesmo tempo têm um número bastante reduzido de linhas de código, o que o torna bastante apelativo.

Em suma, o *Beautiful Soup* revelou ser o módulo ideal para quem quer começar a fazer mineração de dados através de *Web Scraping*, em comparação com outras ferramentas semelhantes e com o mesmo intuito.