

Desenvolvimento Rápido em Linux

Trabalho Prático 4 - Banco de Dados

Webcrawler com Python

Bruno Tomé¹

¹Instituto Federal de Minas Gerais - Campus Formiga (IFMG)
CEP 35570-000 – Formiga – MG – Brasil

`ibrunotome@gmail.com`

Resumo. *Desenvolvimento de um Webcrawler com Python e banco de dados MySQL.*

1. Apresentação

O trabalho consiste no desenvolvimento de um *Webcrawler* para salvar nome, endereço, email e telefone das pessoas cadastradas no site <http://www9.prefeitura.sp.gov.br/secretarias/smads/estouaqui/pessoas/todos/> utilizando a linguagem Python e banco de dados MySQL.

2. Estudo

Foi utilizado a versão 2.7.10 do Python. As bibliotecas utilizadas para a criação do *webcrawler* foram as 4 seguintes:

- *MySQLdb* [1]. *MySQLdb is an thread-compatible interface to the popular MySQL database server that provides the Python database API.*
- *urllib2* [2]. *The urllib2 module defines functions and classes which help in opening URLs (mostly HTTP) in a complex world — basic and digest authentication, redirections, cookies and more.*
- *BeautifulSoup* [3]. *Is a Python library for pulling data out of HTML and XML files. It works with your favorite parser to provide idiomatic ways of navigating, searching, and modifying the parse tree. It commonly saves programmers hours or days of work.*
- *sys* [4]. *This module provides access to some variables used or maintained by the interpreter and to functions that interact strongly with the interpreter. It is always available.*

Os conhecimentos necessários para a realização deste trabalho vieram de um curso chamado “Webcrawler com Python”, da *School of Net* [5], e pode ser acessado neste link: <https://www.schoolofnet.com/curso-webcrawler-com-python/>

3. Caso de uso

O código possui 125 linhas com os comentários. Em uma visão de alto nível, os seguintes passos são realizados:

- 1. Definir link base.
- 2. Ler html.

- 3. Utilizar o soup para ler exatamente as *tags* necessárias para obter informações úteis.
- 4. Criar uma conexão com o banco de dados.
- 5. Inserir os dados enquanto são coletados.
- 6. Repetir o laço até que todas os dados das 10 páginas sejam salvos.

Ao final, temos nome, telefone, email endereço e distrito de 100 pessoas salvos numa tabela “people”.

Abaixo está listado o código para a realização deste trabalho.

```

1 # coding=utf-8
3 # Module for work with mysql database
import MySQLdb
5
6 # The urllib2 module defines functions and classes which help in opening URLs (mostly HTTP)
7 # in a complex world a basic and digest authentication , redirections , cookies and more.
import urllib2
9
10 # Import BeautifulSoup for parse html data
11 from bs4 import BeautifulSoup
13
14 # Import time to know time spent
import time
15
16 import sys
17
18 reload(sys)
19
20 # Link base to make the webcrawler
21 link_base = 'http://www9.prefeitura.sp.gov.br'
link_page = link_base + '/secretarias/smads/estouaqui/pessoas/todos/page:'
23
24 # Headers for browser
25 headers = {'User-Agent': 'GoogleBot'}
27
28 def run_request(link):
29     """
30     Iterate the next 10 pages to get the links and call the function to
31     get the data child
32
33     :param link:
34     :return:
35     """
36
37     global link_page
38
39     if link is None:
40         return
41
42     for i in xrange(1, 11):
43         soup = get_html(link_page + str(i))
44
45         if soup is None:

```

```

47         return
48
49     people_list = soup.find( attrs={ 'id': ' lista_pessoas ' })
50
51     for people in people_list . findAll ( attrs={ 'class': 'pessoa' }):
52         if people is None:
53             continue
54
55         people_data = people . find ( attrs={ 'class': 'pessoa_dados' })
56
57         if people_data is None:
58             return
59
60         link = people_data . h3 . a . get ( ' href ' )
61
62         get_data_child ( link_base + link )
63
64     def get_data_child ( link ):
65         """
66         Get the data of link
67
68         :param link:
69         :return:
70         """
71
72         soup = get_html( link )
73
74         if soup is None:
75             return
76
77         data = soup . find ( attrs={ 'id': 'conteudo' })
78         name = data . h2 . text
79         address = ''
80         district = ''
81         phone = ''
82         email = ''
83
84         count = 1
85
86         # Conection with database
87         conn = connect_db()
88         conn.autocommit(False)
89         cur = conn.cursor ()
90
91         # String for insert into database webcrawler_python, table people, the data crawled from
92         # the link
93         insert = "INSERT INTO people (name, address, district , phone, email) VALUES ('%s', '%s'
94         ', '%s', '%s', '%s')"
95
96         for p in data . find ( attrs={ 'class': ' entidade_dados' }). findAll ( 'p' ):
97
98             if count == 1:
99                 address = p . text . split ( ':' ) [1][1:]
100             elif count == 2:
101                 district = p . text . split ( ':' ) [1][1:]

```

```

101         elif count == 3:
            phone = p.text . split ( ':' ) [1][1:]
103         elif count == 5:
            email = p.text . split ( ':' ) [1][1:]

105         count += 1

107         # Bind the parameters and make the insert
            insert %= (name, address, district , phone, email)
109         cur.execute( insert )
            conn.commit()

111         print name
113
115 def connect_db():
    """
117     Settings for connect to the database
    : return :
119     """

121     return MySQLdb.connect(host='localhost', user='root', passwd='', db='webcrawler_python')
123
125 def get_html( link ):
    """
127     Make the web request for parsing data
    :param link :
    : return :
129     """

131     request = urllib2 .Request( link , headers=headers)
    return BeautifulSoup( urllib2 .urlopen( request ), 'html.parser' )
133
135 if __name__ == '__main__':
    begin = time.time()
137     run_request( link_base + '/ secretarias /smads/estouaqui/pessoas/todos' )
    end = time.time()
139     print '\nTotal: ', (end - begin)

```

Listing 1. Código fonte do trabalho

4. Impressões pessoais

Este minicurso tem uma duração de 6 horas e com ele foi possível ver o quão fácil é tratar html e trabalhar com banco de dados em Python.

Para a utilização do bando de dados MySQL foi necessário apenas importar o MySQLdb, no caso de SQLite seria mais fácil ainda. Já com Django daria um pouco mais de trabalho, porém o Django oferece uma interface de administração pronta, o que é bastante útil.

5. Referências

[1] MySQLdb. Disponível em <<http://mysql-python.sourceforge.net/MySQLdb.html>>

- [2] urllib2. Disponível em <<https://docs.python.org/2/library/urllib2.html>>
- [3] BeautifulSoup. Disponível em <<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>>
- [4] Sys. Disponível em <<https://docs.python.org/2/library/sys.html>>
- [5] School Of Net, Webcrawler com Python. Disponível em <<https://www.schoolofnet.com/curso-webcrawler-com-python/>>