

Образовательное частное учреждение высшего образования
«Еврейский университет»
Факультет экономики и информатики
Кафедра информатики и математики

ЛАБОРАТОРНАЯ РАБОТА № 4
по дисциплине «программирование на Python»

Выполнил: Высоцкий Р.Н.,
студент 3 курса
03.09.09 Прикладная информатика

Руководитель: доцент кафедры информатики
и математики Демичев Василий Анатольевич,
к. ф.-м. н.

Москва, 2022

Цель: получение практических навыков в использовании регулярных выражений и библиотеке BeautifulSoup.

Код:

```
1  import requests
2  from bs4 import BeautifulSoup
3  import re
4
5  import urllib.parse
6
7  def MakeRegFromKeyWords(rawkey):
8      rawkey = rawkey.lower()
9      key = rawkey.split(" ")
10
11     regexp = r"\s";
12     for word in key:
13         regexp += word + r"[a-яA-я]*\s+"
14
15     return regexp
16
```

```
16
17 def ProcessPage(pageurl , depth , regexp , vizited = []):
18     lst = []
19     if depth <= 0:
20         return []
21     pageurl = urllib.parse.unquote(pageurl) #удаление служебных элементов из ссылки
22     if pageurl.find("http") < 0:
23         pageurl = r"http://" + pageurl
24     pageurl = pageurl.replace("index.html" , "")
25     pageurl = pageurl.rstrip("/")
26
27     if pageurl in vizited:
28         return []
29     else:
30         vizited.append(pageurl)
31         print("-->> обрабатываю страницу " , pageurl)
32         try:
33             resp = requests.get(pageurl)
34             soup = BeautifulSoup(resp.text , 'lxml')
35             txt = soup.text
36             txt = txt.lower()
37         except:
38             print("---->> ссылка " , pageurl , "не работает!")
39             return []
40     match = re.search(regexp , txt)
41     if match:
42         start = match.start()
43         end = match.end()
44         print("-->> найдено совпадение на " , pageurl)
45
46         s0 = start - 100
47         if s0 < 0:
48             s0 = 0
49         e0 = end + 100
50         if e0 >= len(txt):
51             e0 = len(txt)-1
52
53         print(txt[s0:e0])
54         lst.append([pageurl , start , end , txt[s0:e0]])
```

```

55
56     for tag in soup.find_all("a" , href = True):
57         mas = ProcessPage(tag['href'] , depth-1 , regexp , vizited)
58         for items in mas:
59             lst.append(items)
60
61     return lst
62
63 address = input("Ведите адрес сайта: ")
64 key = input("Введите ключевую фразу для поиска: ")
65 depth = int(input("Введите поисковую глубину: "))
66 regexp = MakeRegFromKeyWords(key)
67 vizited = []
68 mas = ProcessPage(address , depth , regexp , vizited)
69 percent = len(mas)/len(vizited) * 100.0
70 print("> процент совпадения ключевой фразы: " , percent)
71 regexp = r"https?:?\/\/(www\.)?.\w+\.\w+(\.\w+)?(\.\w+)?(\.\w+)?(\.\w+)?"
72 #убираем название сайта
73 websites = {}
74 for item in mas:
75     match = re.search(regexp , item[0])
76     if match:
77         start = match.start()
78         end = match.end()
79         website = item[0][start:end]
80         if website in websites.keys():
81             websites[website] += 1
82         else:
83             websites[website] = 1
84
85 print("наличие информации на сайтах:" , websites)
86
87

```

Программа рекурсивно ищет слова методом из предыдущей лабораторной работы, а так же через BeautifulSoup ищет ссылки на странице и переходит по ним для дальнейшего поиска.