

目录

1	实验介绍	2
1.1	实验内容	2
1.1.1	实验二	2
1.1.2	实验三	2
1.2	实验环境	2
1.3	所需技能	2
2	实验过程	3
2.1	Basic Crawler	3
2.1.1	BBS 修改个人说明档	3
2.1.2	BFS	3
2.1.3	爬取与解析页面	4
2.1.4	网页编码自动检测	4
2.2	Bloom Filter	5
2.2.1	Hash Function	5
2.2.2	Bit Array	6
2.2.3	BloomFilter 类实现	7
2.3	Multi-threading	8
3	总结与分析	10

1 实验介绍

1.1 实验内容

本次报告涵盖了第二次和第三次实验, 其内容分别如下:

1.1.1 实验二

1. 使用自己的账号模拟登陆 BBS 后, 修改个人说明档 (修改 `bbs_set_sample.py`).
2. 修改 `crawler_sample.py` 中的 `union_bfs` 函数, 完成 BFS 搜索.
3. 修改 `crawler_sample.py` 中的 `crawl` 函数, 返回图的结构.
4. 进一步修改函数, 完成网页爬虫 (将练习 2,3 中修改的部分加入 `crawler.py`, 注意做异常处理).

1.1.2 实验三

1. 实现 BloomFilter, 并设计一个实验统计你的 BloomFilter 的错误率 (false positive rate).
2. 实现一个并行的爬虫, 将实验二中的 `crawler.py` 改为并行化实现.

1.2 实验环境

操作系统: Ubuntu 18.04 LTS

Python 版本: 3.7

1.3 所需技能

实验二主要涉及运用 Python 模拟网络请求, 这需要具备一定的网络知识. 而其后对网页内容进行解析, 并根据不同策略进行爬取, 则需要基本的算法知识.

实验三涉及到 Hash Function, Hash Table, BitMap 以及 Bloom Filter 的原理以及实现. 在编写多线程并行程序时, 还需要对操作系统, 阻塞以及锁等概念有些了解.

2 实验过程

2.1 Basic Crawler

2.1.1 BBS 修改个人说明档

该练习主要在于了解用 POST 请求提交数据的流程, 较为简单.

实验提供的指导是用 urllib 库来模拟网络请求, cookielib 库来操作 cookie, 操作较为繁琐. 事实上, 由于本练习仅需要保持登录状态, 并不需要涉及直接操作 cookie 的细节, 因此可以使用 requests 库提供的 Session() 对象. 如下:

```
1 def bbs_set(id, pw, text):
2
3     s = requests.Session()
4     login_data = {
5         "id": id,
6         "pw": pw,
7         "submit": "login",
8     }
9     login_url = "https://bbs.sjtu.edu.cn/bbslogin"
10    login_request = s.post(login_url, login_data)
11    assert login_request.status_code == 200, "Failed to log in."
12    update_data = {
13        "text": text,
14        "type": "update",
15    }
16    update_url = "https://bbs.sjtu.edu.cn/bbsplan"
17    update_request = s.post(update_url, update_data)
18    content = s.get('https://bbs.sjtu.edu.cn/bbsplan').content
19    soup = BeautifulSoup(content, features="html.parser")
20    print(str(soup.find('textarea').string).strip())
```

值得注意的是对提交本文的编码, 需要进行 GBK 编码.

2.1.2 BFS

由于此处未考虑运行的效率, BFS 函数的实现相当简单:

```
1 def union_bfs(a, b):
2     for e in b:
3         if e not in a:
4             a.insert(0, e)
```

2.1.3 爬取与解析页面

在爬取页面时, 为了避免程序中途因为网络原因而崩溃, 需要捕获异常以增强健壮性:

```
1 def get_page(page):
2     try:
3         r = requests.get(page, timeout=10)
4     except:
5         return None
6     content = r.content
7     return content
```

而在解析页面时, 可以复用实验一时所写的 `BaseParser` 类, 而无需重复编码:

```
1 from parsers import BaseParser
2 def get_all_links(content, page):
3     parser = BaseParser()
4     url_set = parser.parse_url(content, page)
5     links = [u for u in url_set]
6     return links
```

2.1.4 网页编码自动检测

在实验一的报告中, 我将网页编码的问题暂时搁置. 而在这次实验中, 由于要爬取的目标不再是单一的网站, 而是互联网上各式各样的网站, 其编码方式各异, 网页编码的问题便无法回避了.

`add_page_to_folder` 函数将爬取的网页保存到本地文件 (UTF-8 编码), 其间涉及到两次编码解码: 一是对网页内容解码, 二是在写入文件时编码. 第二步并无大碍, 但第一步, 如何对不同编码的网页解码呢?

事实上, Python 提供了 `chardet` 库专门用以检测文本的编码, 只需在解码前先检测出可能的编码, 便可以从容应对大多数的情况. 至于少数检测不出编码方式的网页文本, 只能按照默认的 UTF-8 来处理, 并忽略解码时的错误. 如下:

```
1 import chardet
2
3 def add_page_to_folder(page, content):
4     #some code omitted
5     #...
6     f = open(os.path.join(folder, filename), 'w', encoding="utf8")
7     encode = chardet.detect(content)
8     if not encode["encoding"]:
9         print("Failed to detect encoding.({url},{encode})".format(url=page,
10             ↪ encode=encode))
11     f.write(content.decode(encode["encoding"] or "utf8", errors="ignore"))
12     f.close()
```

2.2 Bloom Filter

2.2.1 Hash Function

网络上有许多种公开可用的哈希函数, 像本次实验提供的 General Purpose Hash Function Algorithms Library 中就已有 10 余种哈希函数, 除此之外, 更是有诸如 MD5, SHA1 等常用于加密的哈希函数.

对于 Bloom Filter, 哈希函数的选用十分关键, 关涉到运行的效率与检测的正确性. 考虑到我们的过滤器仅仅用于爬虫的网址去重, 相较于安全性, 更注重运行的速度.

那么, 如何找到一个高效率的哈希函数呢? 诚然, 可以从算法角度进行分析, 但既然我们手中已有现成的代码实现, 设计一个简单进行测试是最直接的做法. 为此, 我编写了 HashTest.py, 用以比较不同的哈希函数的运行效率:

```
1 from GeneralHashFunctions import *
2 import mmh3
3 import random
4 import time
5 from string import ascii_letters as letters

6 def test_hash_efficiency(hash_function, test_strings):
7     start_time = time.time()
8     for s in test_strings:
9         hash_function(s)
10    finish_time = time.time()
11    return finish_time - start_time

12 test_strings = ["".join(random.sample(letters, 32)) for i in range(100000)]
13 hash_functions = [hash, mmh3.hash, JSHash, PJWHash, ELFHash, BKDRHash,
14 ↪ SDBMHash, DJBHash, DEKHash, BPHash, FNVHash, APhash]
15 result = []

16 for f in hash_functions:
17     result.append((f.__name__, test_hash_efficiency(f, test_strings)))
18     print("{0:<9}: {1:.8f}s".format(*result[-1]))
19 print("*****Result Sorted By Time*****")
20 result.sort(key=lambda x: x[1])
21 for r in result:
22     print("{0:<9}: {1:.8f}s({2:.3f})x".format(*r, r[1]/result[0][1]))
```

随机生成了 100000 个长度为 32 的字符串, 然后对其用不同的哈希函数进行散列处理, 比较运行时间, 最终得到的结果如下:

```

1 *****Result Sorted By Time*****
2 hash      : 0.00897479s(1.000x)  # Built-in Hash
3 hash      : 0.02493358s(2.778x)  # mmh3.hash
4 BPHash    : 0.73630977s(82.042x)
5 BKDRHash  : 0.83501196s(93.040x)
6 DJBHash   : 0.88225102s(98.303x)
7 FNVHash   : 0.99035144s(110.348x)
8 DEKHash   : 1.01107359s(112.657x)
9 SDBMHash  : 1.23072314s(137.131x)
10 PJWHash   : 1.63085604s(181.715x)
11 JSHash    : 2.06598854s(230.199x)
12 APhash    : 2.06879592s(230.512x)
13 ELFHash   : 2.11861658s(236.063x)

```

可见,从运行效率上来说,不同哈希函数的差距巨大(相差最大236倍),其中,Python自带的`hash()`最优, `murmurhash3`紧随其后.但是,考虑到在 Bloom Filter 中需要对同一个哈希函数加以不同的种子(当然也可以用不同的哈希函数,但那样较为繁琐),而 Python 自带的`hash()`并不支持这么做,故最终我选择了 `murmurhash3`, 近年兴起的一种高效率哈希函数.

2.2.2 Bit Array

在实验提供的 Python 中,已实现了 `Bitarray` 类,但由于其编写于 Python 2.x 环境下,并且存在一个严重的 Bug,需要对其进行一些修正.修改后的 `Bitarray` 类如下:

```

1 class Bitarray:
2     def __init__(self, size):
3         """ Create a bit array of a specific size """
4         self.size = size
5         self.bitarray = bytearray(size//8 + 1) # Attention Here!!!
6
7     def set(self, n):
8         """ Sets the nth element of the bitarray """
9         index = n // 8
10        position = n % 8
11        self.bitarray[index] = self.bitarray[index] | 1 << (7 - position)
12
13    def get(self, n):
14        """ Gets the nth element of the bitarray """
15        index = n // 8
16        position = n % 8
17        return (self.bitarray[index] & (1 << (7 - position))) > 0

```

需要注意的是,原程序对 `bytearray` 的大小设置有误,运行时有概率导致下标越界.

2.2.3 BloomFilter 类实现

对于哈希函数个数 k 、位数组大小 m 、加入的字符串数量 n , 根据一定的数学推导, 可知当 $k = \ln(2) * m/n$ 时出错的概率是最小的. 而设定 $m = 20 * n$, 则出错的概率为 0.0000889, 对于爬虫而言完全可以接受.

因此, 在类初始化时, 既允许同时指定 n, m, k , 也可以允许只指定 n , 其余值按照上述规则计算.

而在对哈希函数列表初始化时, 用到了 Python 具有的函数式语言特点, 通过自行实现的类的私有方法 `__hash_function_wrapper` 来对哈希函数进行处理, 方便之后的调用:

```
1 import mmh3      # third-party library for murmurhash3
2 from functools import partial
3 import math
4
5 class BloomFilter(object):
6
7     def __hash_function_wrapper(self, hash_func, seed, mod):
8         def f(s):
9             return partial(hash_func, seed=seed)(s) % mod
10        return f
11
12    def __init__(self, potential_num, bitarray_size=None, hash_num=None):
13        self.potential_num = potential_num
14        self.bitarray_size = bitarray_size or 20 * potential_num
15        self.hash_num = hash_num or
16        ↪ int(math.log(2)*self.bitarray_size/self.potential_num)
17        self.hash_functions = [self.__hash_function_wrapper(
18            mmh3.hash, seed=i, mod=self.bitarray_size
19        ) for i in range(self.hash_num)]
20        self.bitarray = Bitarray(self.bitarray_size)
21        self.count = 0
```

在初始化完成之后, Bloom Filter 所需实现的操作为加入元素 `add` 与检查元素是否已加入 `check`, 至于 `__len__` 与 `parameters`, 是为了方便外界获知 Bloom Filter 的状态.

```
1 def add(self, string):
2     for hash_func in self.hash_functions:
3         self.bitarray.set(hash_func(string))
4     self.count += 1
5
6 def check(self, string):
7     return all(self.bitarray.get(hash_func(string))
8                for hash_func in self.hash_functions)
```

```

8 def __len__(self):
9     return self.count

10 def parameters(self):
11     return self.potential_num, self.bitarray_size, self.hash_num

```

2.3 Multi-threading

在原有的 crawler.py 程序基础上, 为了使其能够以多线程的方式并发爬取, 需要修改 `crawl` 函数.

为了不引入过多的全局变量污染程序, 利用 Python 的函数式特性: 闭包, 我在 `crawl` 函数内部定义了 `crawl_single_page` 函数作为多线程的 worker. 大致如下:

```

1 def crawl(seed, max_page, thread_num=4):

2     def crawl_single_page():
3         pass
4     url_queue = Queue()
5     lock = Lock()
6     crawled = BloomFilter(max_page)
7     graph = {}
8     global count
9     count = 0
10    url_queue.put(seed)
11    for i in range(thread_num):
12        t = Thread(target=crawl_single_page)
13        t.setDaemon(True)
14        t.start()
15    url_queue.join()
16    return graph, crawled

```

需要注意的是, 不同于实验二, 这里的 `crawled` 不是列表, 而是我们所实现的 Bloom Filter.

至于 `crawl_single_page` 的具体实现, 与实验二基本无异, 不同之处在于锁的引入, 避免多线程导致变量读写的异常.

```

1 def crawl_single_page():
2     global count
3     while count < max_page:
4         url = url_queue.get()
5         if not crawled.check(url):
6             print("#{0:<4} {1}".format(count+1, url))
7             content = get_page(url)

```



```

8         outlinks = get_all_links(content, url)
9         with lock:
10             if count >= max_page: #attention here
11                 print("Drop!")
12                 break
13             graph[url] = outlinks
14             for l in outlinks:
15                 url_queue.put(l)
16             crawled.add(url)
17             count += 1
18         add_page_to_folder(url, content)
19     url_queue.task_done()

```

留意以上代码, 为了确保最终爬取的页面数量与设定一致, 在爬取完页面后, 需要丢弃由于多线程导致的冗余页面 (如目标是 100 个页面, 但在爬取完 99 个页面后, 4 个线程从队列中取出一个页面进行爬取).

然而, 由于我们在爬取到足够数量页面之后就停止爬取, 队列中往往会存留尚未爬取的页面, 这就导致了在线程结束后依旧阻塞, 程序无法终止. 因此, 在主程序中, 应对每个 Thread 进行 join 操作, 而非对 Queue.

```

1 threads = []
2 for i in range(thread_num):
3     t = Thread(target=crawl_single_page)
4     threads.append(t)
5     t.setDaemon(True)
6     t.start()
7 for t in threads:
8     t.join()

```

3 总结与分析

经过实验二与实验三,相较于实验一中简单的爬虫,已有不少提升,如 **Bloom Filter** 和多线程对运行效率的显著提升,以及许多诸如网页编码自动检测等细节的优化.在进行实验的过程中,课程提供的指引起到了很大的作用,但在某些情形下,还是需要自己查找资料并动手实践来解决问题,这样让我学到了不少新的内容.

此外,每次实现一个功能,总是有不止一种可供选择的方案.有的方案简单而粗暴,却缺乏美感(譬如引入一大堆全局变量),我总是试图避免出现这样的情况,因而或许写出了一些不太常规的代码.

虽然多费了时间,但挺有趣的.