# Report of Lab3

Yuxiang Chen 5110309783

October 21, 2012

## Contents

# 1 The Purpose of Lab3 and My Preparation

It not hard to find that the purpose of lab3 is to let us know and command the theorem of indexing and searching the information on the Internet based on the use of lucene. And what we need to know first is the principle of Full-text Search. In this method, we need to firstly index the information we have got on the Internet, then search the index according to some key words. Although the principle sounds easy, it's a little hard for me to do it very fast, since I don't know what lucene is and how it works at first. Then I have to study on my own ...... And at last I find how I can accomplish my experiment.

The first part of the Full-text Search is indexing, and it is composed of the following 4 parts:

1. Find the document you want to index;
2. Pass the document into the Tokenizer;
3. Give the token you got at step 2 to the linguistic processor;
4. Pass the term to the indexer.

And the latter half – search can also be split into 4 parts:

1. The user input the query;
2. Parse and process the query;
3. Search the index and find the document agreed to the syntax tree;
4. Sort the result according to the relativity of the document, etc.

The second half of the experiment is mainly an extending on the basis of the first part, which acquires us to contain more information indexed and searched in our programs. The first needs an extra information 'site' and the second one needs the url of the picture to be searched, and the name and url of the page it belongs to.

And thus using the tools we have, it's time to achieve our goal.

# 2 The Main Part of the Experiment

In this experiment, I split the program into three parts: crawling to get web pages, setting up index, and searching according to the query.

## 2.1 The First Half of the Experiment

The first part is separated into three parts by me, which in my view makes the experiment more organized and easier to understand.

### 2.1.1 Crawl to Get Web Pages

In this part, I mainly used the program we wrote in lab2 to crawl web pages and store them in "F:\html". But I found that the speed to crawl web pages is a little slow, so I choose 'http://www.sjtu.edu.cn' to make it fast for us to get enough pages. If you want to crawl other pages, you may just change it in the following code.

And here is the code:

```
from BeautifulSoup import BeautifulSoup
import urllib2
import re
import urlparse
import os
import urllib
import socket
import threading
import Queue
import time
import chardet
import sys
reload(sys)
```

```python
sys.setdefaultencoding('utf8')


def valid_filename(s):
    import string
    valid_chars = "-_.() _%s%s" % (string.ascii_letters, string.digits)
    s = ''.join(c for c in s if c in valid_chars)
    return s

def get_page(page):
    time.sleep(0.001)
    try:
        content=urllib2.urlopen(page,timeout=3).read()
        result = chardet.detect(content)['encoding']
        if result=='GB2312':
            content=content.decode('gbk').encode('utf8')
        return content
    except:
        #There is an error.#
        return []

def get_all_links(content, page):
    if content==[]:
        return []
    links = []
    tempset=set()
    soup=BeautifulSoup(content)
    for i in soup.findAll('a',{'href':re.compile(('^http|^/'))}):
        tempset.add(i['href'])
    for i in tempset:
        links.append(urlparse.urljoin(page,i))
    return links

def add_page_to_folder(page, content):
    folder = 'F:\\html'
    index_filename = 'F:\\html\index.txt'
    filename = valid_filename(page)
    index = open(index_filename, 'a')
    index.write(filename + ';' + page + '\n')
    index.close()
    if not os.path.exists(folder):
        os.mkdir(folder)
    f = open(os.path.join(folder, filename), 'w')
    f.write(content)
    f.close()
```

```
60  def working ():
61      page_num=0
62      while page_num<task_per_thread:
63          page = q.get()
64          if page not in crawled:
65              content = get_page(page)
66              outlinks = get_all_links(content,page)
67              if outlinks==[]:
68                  q.task_done()
69                  continue
70              page_num+=1
71              add_page_to_folder(page,content)
72              for link in outlinks:
73                  q.put(link)
74              if varLock.acquire():
75                  crawled.append(page)
76                  varLock.release()
77                  q.task_done()
78              else:
79                  q.task_done()
80          else:
81              q.task_done()
82      while q.empty()==False:
83          q.get()
84          q.task_done()
85
86  if not os.path.exists("F:\\html"):
87      os.mkdir("F:\\html")
88  NUM = 100
89  task_per_thread=50
90  crawled = []
91  varLock = threading.Lock()
92  q = Queue.Queue()
93  q.put('http://www.sjtu.edu.cn')
94  for i in range(NUM):
95      t=threading.Thread(target=working)
96      t.setDaemon(True)
97      t.start()
98  q.join()
99  print "That's all you want."
```

In this part, we need to analyze the coding method of the page, and I transform them all to 'utf8', which will make the following parts easier. And then I save the name of the page and its url to "F:\html\index.txt". And here are the pictures of the compiler and the file. And in the picture, plus one "index.txt", there are 5001 documents in total.

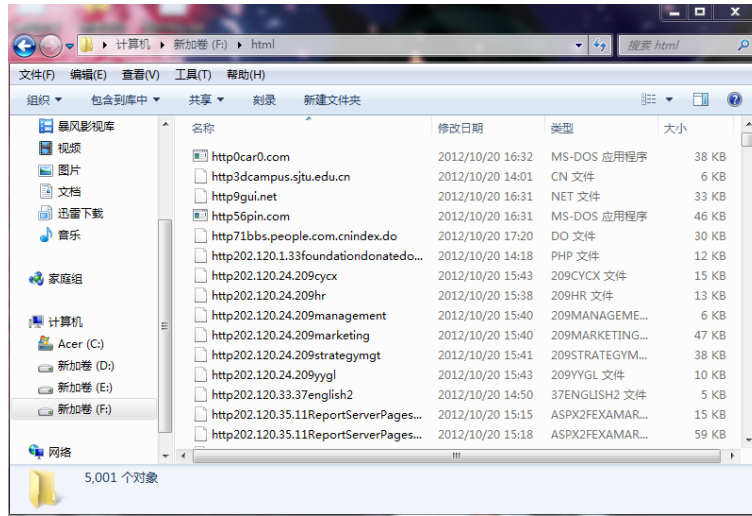Figure 1: the outcome of the compiler



Figure 2: the screenshot of F:\html

### 2.1.2 Set up the Index

It is really a little annoying in this part, since there're so many things to do. In this part, we have to set up an index to save the name,path,title,url and contents of the page we have crawled. Meanwhile, I also give the Chinese words a segmentation using ICTCLAS50. And though it's quite complex, I feel it is hard to describe how to achieve them, so I will only display the code and you may find how it works. By the way, in order to find url easily, I make a dictionary using 'index.txt' we saved in the first part.

```
1  import sys, os, lucene, threading, time, chardet, urllib2
2  from datetime import datetime
3  from BeautifulSoup import BeautifulSoup
4  from ctypes import *
5
```

Figure 3: the screenshot of index.txt

```
 6  """
 7  This class is loosely based on the Lucene (java implementation) demo class
 8  org.apache.lucene.demo.IndexFiles.  It will take a directory as an argument
 9  and will index all of the files in that directory and downward recursively.
10  It will index on the file path, the file name and the file contents.
    The
11  resulting Lucene index will be placed in the current directory and called
12  'index'.
13  """
14
15  class IndexFiles(object):
16      """Usage: python IndexFiles <doc_directory>"""
17
18      def __init__(self, root, storeDir, analyzer):
19
20          if not os.path.exists(storeDir):
21              os.mkdir(storeDir)
22          store = lucene.SimpleFSDirectory(lucene.File(storeDir))
23          writer = lucene.IndexWriter(store, analyzer, True,
24                                       lucene.IndexWriter.MaxFieldLength.LIMITED)
25          writer.setMaxFieldLength(1048576)
26          self.indexDocs(root, writer)
27          print 'optimizing index.',
28          writer.optimize()
29          writer.close()
30          print 'done'
```

```python
     def indexDocs(self, root, writer):
         for root, dirnames, filenames in os.walk(root):
             for filename in filenames:
                 if filename.endswith('.txt'):
                     continue
                 print "adding", filename
                 try:
                     path = os.path.join(root, filename)
                     file = open(path)
                     buf = file.read()
                     contents=buf
                     result = chardet.detect(buf)['encoding']
                     if result=='GB2312':
                         contents = buf.decode('gbk').encode('utf8')
                     file.close()
                     soup=BeautifulSoup(contents)
                     url=mydict[filename]
                     title=str(soup.head.title.string).decode('utf8')
                     contents=''.join(soup.findAll(text=True))
                     doc = lucene.Document()
                     doc.add(lucene.Field("name", filename,
                                          lucene.Field.Store.YES,
                                          lucene.Field.Index.NOT_ANALYZED))
                     doc.add(lucene.Field("path", path,
                                          lucene.Field.Store.YES,
                                          lucene.Field.Index.NOT_ANALYZED))
                     doc.add(lucene.Field("url", url,
                                          lucene.Field.Store.YES,
                                          lucene.Field.Index.NOT_ANALYZED))
                     doc.add(lucene.Field("title", title,
                                          lucene.Field.Store.YES,
                                          lucene.Field.Index.NOT_ANALYZED))
                     if len(contents) > 0:
                         dll=cdll.LoadLibrary("F:\\ICTCLAS50_Windows_32_C\ICTCLAS
                         dll.ICTCLAS_Init(c_char_p("F:\\ICTCLAS50_Windows_32_C"))
                         strlen = len(c_char_p(contents).value)
                         t =c_buffer(strlen*6)
                         bSuccess = dll.ICTCLAS_ParagraphProcess
                         (c_char_p(contents),c_int(strlen),t,c_int(0),0)
                         contents=t.value.decode('gbk').encode('utf8')
                         ##list=t.value.split()
                         ##print ' '.join(list)
                         dll.ICTCLAS_Exit()
                         doc.add(lucene.Field("contents", contents,
                                              lucene.Field.Store.NO,
```

```
77                                                    lucene.Field.Index.ANALYZED))
78                          else:
79                              print "warning:_no_content_in_%s" % filename
80                          writer.addDocument(doc)
81                    except Exception, e:
82                        print "Failed_in_indexDocs:", e
83
84   if __name__ == '__main__':
85   ##      if len(sys.argv) < 2:
86   ##          print IndexFiles.__doc__
87   ##          sys.exit(1)
88        lucene.initVM()
89        print 'lucene', lucene.VERSION
90        start = datetime.now()
91        dic= open('F:\\html\index.txt')
92        d = dic.readlines()
93        dic.close()
94        mydict = {}
95        for word in d:
96            key = word.split(';')[0]
97            value = word.split(';')[1]
98            mydict[key] = value
99        try:
100  ##          IndexFiles(sys.argv[1], "index", lucene.SimpleAnalyzer(lucene.Version..
101             IndexFiles('F:\\html', "F:\\index", lucene.SimpleAnalyzer(lucene.Version
102             end = datetime.now()
103             print end - start
104        except Exception, e:
105            print "Failed:_", e
```

So after this part, we can get the index set up in F:\index, making the last part easily to be done. And I use 'SimpleAnalyzer' in the indexing and searching part, which will make the whole experiment more execute.
And below are the pictures of the outcome in the interpreter and the index in file "F:\index".

### 2.1.3  Search According to the Query

Things get really easy when they come to the last part. In this searching part, we only need to change the example a little bit, adding the outcome of title and url, which we have already made an index in the second part.
And that are the codes:

```
1   from lucene import \
2       QueryParser, IndexSearcher, SimpleAnalyzer, SimpleFSDirectory, File, \
```

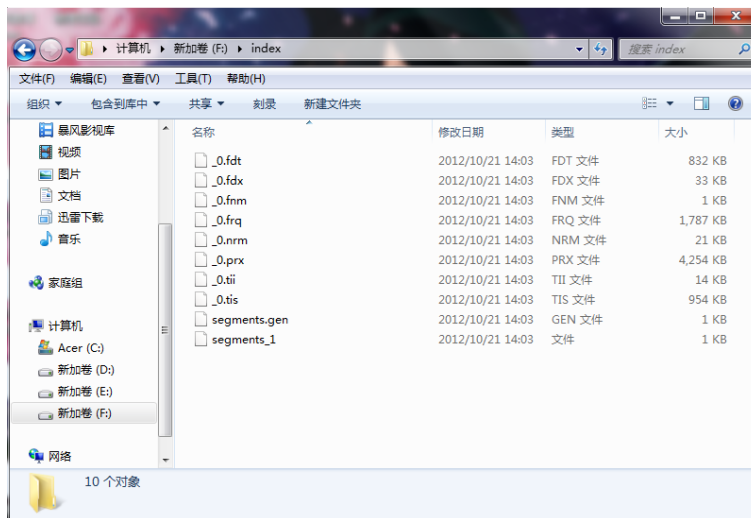Figure 4: the outcome in interpreter



Figure 5: the screenshot of F:\index

```
 3        VERSION, initVM, Version
 4
 5
 6    """
 7    This script is loosely based on the Lucene (java implementation) demo class
 8    org.apache.lucene.demo.SearchFiles.  It will prompt for a search query, then it
 9    will search the Lucene index in the current directory called 'index' for the
10    search query entered against the 'contents' field.  It will then display the
11    'path' and 'name' fields for each of the hits it finds in the index.
      Note that
12    search.close() is currently commented out because it causes a stack overflow in
13    some cases.
14    """
15    def run(searcher, analyzer):
16        while True:
17            print
18            print "Hit_enter_with_no_input_to_quit."
19            command = raw_input("Query:")
20            if command == '':
21                return
22            print
23            print "Searching_for:", command
24            query = QueryParser(Version.LUCENE_CURRENT, "contents",
25                                analyzer).parse(command)
26            scoreDocs = searcher.search(query, 50).scoreDocs
27            print "%s_total_matching_documents." % len(scoreDocs)
28
29            for scoreDoc in scoreDocs:
30                doc = searcher.doc(scoreDoc.doc)
31                print 'path:', doc.get("path"), 'title:', doc.get("title"),'url:', d
32
33
34    if __name__ == '__main__':
35        STORE_DIR = "F:\\index"
36        initVM()
37        print 'lucene', VERSION
38        directory = SimpleFSDirectory(File(STORE_DIR))
39        searcher = IndexSearcher(directory, True)
40        analyzer = SimpleAnalyzer(Version.LUCENE_CURRENT)
41        run(searcher, analyzer)
42        searcher.close()
```

Then we can search the information in the pages we have crawled. These are the screenshots:

And thus, at last I finished the experiment working so hard, but luckily, this time, I was not troubled with the code problem as often as I used to be, thanks

Figure 6: the outcome in interpreter(1)

Figure 7: the outcome in interpreter(2)

Figure 8: the outcome in interpreter(3)

Figure 9: the outcome in interpreter(4)



Figure 10: the outcome in interpreter(5)

to the "Q.ppt".

## 2.2 The Second part of the Experiment

It easy to find that the crawler part of the latter half of lab3 is nothing different from the first part, so I will only give out the index and search part of each experiment.

### 2.2.1 Make Index with the Information of Site and Search It

There is little different between this program and the indexing program of the first part, except that we need to find the domain name of the website and add it to the final index. And I happened to find an easy way to get the domain name using urllib, and I will present it in the following codes. By the way, I use the word 'site' instead of 'domain name' in my program.

```python
import sys, os, lucene, threading, time, chardet, urllib2
from datetime import datetime
from BeautifulSoup import BeautifulSoup
from ctypes import *
import urllib

"""
This class is loosely based on the Lucene (java implementation) demo class
org.apache.lucene.demo.IndexFiles.  It will take a directory as an argument
and will index all of the files in that directory and downward recursively.
It will index on the file path, the file name and the file contents.  The
resulting Lucene index will be placed in the current directory and called
'index'.
"""

class Ticker(object):

    def __init__(self):
        self.tick = True

    def run(self):
        while self.tick:
            sys.stdout.write('.')
            sys.stdout.flush()
            time.sleep(1.0)

class IndexFiles(object):
    """Usage: python IndexFiles <doc_directory>"""

```

```
30      def __init__(self, root, storeDir, analyzer):
31
32          if not os.path.exists(storeDir):
33              os.mkdir(storeDir)
34          store = lucene.SimpleFSDirectory(lucene.File(storeDir))
35          writer = lucene.IndexWriter(store, analyzer, True,
36                                      lucene.IndexWriter.MaxFieldLength.LIMITED)
37          writer.setMaxFieldLength(1048576)
38          self.indexDocs(root, writer)
39          ticker = Ticker()
40          print 'optimizing index',
41          threading.Thread(target=ticker.run).start()
42          writer.optimize()
43          writer.close()
44          ticker.tick = False
45          print 'done'
46
47      def indexDocs(self, root, writer):
48          for root, dirnames, filenames in os.walk(root):
49              for filename in filenames:
50                  if filename.endswith('.txt'):
51                      continue
52                  print "adding", filename
53                  try:
54                      path = os.path.join(root, filename)
55                      file = open(path)
56                      buf = file.read()
57                      contents=buf
58                      result = chardet.detect(buf)['encoding']
59                      if result=='GB2312':
60                          contents = buf.decode('gbk').encode('utf8')
61                      file.close()
62                      soup=BeautifulSoup(contents)
63                      url=mydict[filename]
64                      proto, rest = urllib.splittype(url)
65                      site, rest = urllib.splithost(rest)
66                      title=str(soup.head.title.string).decode('utf8')
67                      contents=''.join(soup.findAll(text=True))
68                      doc = lucene.Document()
69                      doc.add(lucene.Field("name", filename,
70                                           lucene.Field.Store.YES,
71                                           lucene.Field.Index.NOT_ANALYZED))
72                      doc.add(lucene.Field("path", path,
73                                           lucene.Field.Store.YES,
74                                           lucene.Field.Index.NOT_ANALYZED))
75                      doc.add(lucene.Field("url", url,
```

```
76                                         lucene.Field.Store.YES,
77                                         lucene.Field.Index.NOT_ANALYZED))
78                   doc.add(lucene.Field("title", title,
79                                         lucene.Field.Store.YES,
80                                         lucene.Field.Index.NOT_ANALYZED))
81                   doc.add(lucene.Field("site", site,
82                                         lucene.Field.Store.YES,
83                                         lucene.Field.Index.ANALYZED))
84                   if len(contents) > 0:
85                       dll=cdll.LoadLibrary("F:\\ICTCLAS50_Windows_32_C\ICTCLAS
86                       dll.ICTCLAS_Init(c_char_p("F:\\ICTCLAS50_Windows_32_C"))
87                       strlen = len(c_char_p(contents).value)
88                       t =c_buffer(strlen*6)
89                       bSuccess = dll.ICTCLAS_ParagraphProcess
90                       (c_char_p(contents),c_int(strlen),t,c_int(0),0)
91                       contents=t.value.decode('gbk').encode('utf8')
92                       ##list=t.value.split()
93                       ##print ' '.join(list)
94                       dll.ICTCLAS_Exit()
95                       doc.add(lucene.Field("contents", contents,
96                                         lucene.Field.Store.NO,
97                                         lucene.Field.Index.ANALYZED))
98                   else:
99                       print "warning: no content in %s" % filename
100                  writer.addDocument(doc)
101              except Exception, e:
102                  print "Failed in indexDocs:", e
103
104 if __name__ == '__main__':
105 ##       if len(sys.argv) < 2:
106 ##           print IndexFiles.__doc__
107 ##           sys.exit(1)
108     lucene.initVM()
109     print 'lucene', lucene.VERSION
110     start = datetime.now()
111     dic= open('F:\\html\index.txt')
112     d = dic.readlines()
113     dic.close()
114     mydict = {}
115     for word in d:
116         key = word.split(';')[0]
117         value = word.split(';')[1]
118         mydict[key] = value
119     try:
120 ##           IndexFiles(sys.argv[1], "index", lucene.WhitespaceAnalyzer(lucene.Vers
121          IndexFiles('F:\\html', "F:\\index", lucene.WhitespaceAnalyzer(lucene.Ver
```

15

```
122            end = datetime.now()
123            print end − start
124        except Exception, e:
125            print "Failed: ", e
```

As for the search part, we only need to add the keyword 'site' in the program to make it possible to search the contents in a certain site. In this part, we have to use BooleanQuery to accomplish the query consisted of several kinds of keywords, which can be learnt in the ppt file given.

```
 1  from lucene import \
 2      QueryParser, IndexSearcher, WhitespaceAnalyzer, SimpleFSDirectory, File, \
 3      VERSION, initVM, Version, BooleanQuery, BooleanClause
 4
 5
 6  """
 7  This script is loosely based on the Lucene (java implementation) demo class
 8  org.apache.lucene.demo.SearchFiles. It will prompt for a search query, then it
 9  will search the Lucene index in the current directory called 'index' for the
10  search query entered against the 'contents' field. It will then display the
11  'path' and 'name' fields for each of the hits it finds in the index. Note that
12  search.close() is currently commented out because it causes a stack overflow in
13  some cases.
14  """
15
16  def parseCommand(command):
17      '''
18      input: C title:T author:A language:L
19      output: {'contents':C, 'title':T, 'author':A, 'language':L}
20
21      Sample:
22      input:'contenance title:henri language:french author:william shakespeare'
23      output:{'author': ' william shakespeare',
24              'language': ' french',
25              'contents': ' contenance',
26              'title': ' henri'}
27      '''
28      allowed_opt = ['site']
29      command_dict = {}
30      opt = 'contents'
31      for i in command.split(' '):
32          if ':' in i:
33              opt, value = i.split(':')[:2]
34              opt = opt.lower()
35              if opt in allowed_opt and value != '':
36                  command_dict[opt] = command_dict.get(opt, '') + ' ' + value
```

```python
37              else:
38                  command_dict[opt] = command_dict.get(opt, '') + ' ' + i
39      return command_dict
40
41
42  def run(searcher, analyzer):
43      while True:
44          print
45          print "Hit enter with no input to quit."
46          command = raw_input("Query:")
47          if command == '':
48              return
49
50          print
51          print "Searching for:", command
52
53          command_dict = parseCommand(command)
54          querys = BooleanQuery()
55          for k,v in command_dict.iteritems():
56              query = QueryParser(Version.LUCENE_CURRENT, k,
57                                  analyzer).parse(v)
58              querys.add(query, BooleanClause.Occur.MUST)
59          scoreDocs = searcher.search(querys, 50).scoreDocs
60          print "%s total matching documents." % len(scoreDocs)
61
62          for scoreDoc in scoreDocs:
63              doc = searcher.doc(scoreDoc.doc)
64  ##            explanation = searcher.explain(query, scoreDoc.doc)
65              print "————————————————————"
66              print 'path:', doc.get("path")
67              print 'name:', doc.get("name")
68              print 'title:', doc.get('title')
69              print 'url:', doc.get('url')
70  ##             print explanation
71
72
73  if __name__ == '__main__':
74      STORE_DIR = "F:\\index"
75      initVM()
76      print 'lucene', VERSION
77      directory = SimpleFSDirectory(File(STORE_DIR))
78      searcher = IndexSearcher(directory, True)
79      analyzer = WhitespaceAnalyzer(Version.LUCENE_CURRENT)
80      run(searcher, analyzer)
81      searcher.close()
```

In this part, to make it fast, I only get 30 pages as an example, you can crawl more pages by modifying the variable in the crawler program. And here are some screenshots of the effect of my program.



Figure 11: the outcome of the site part(1)



Figure 12: the outcome of the site part(2)

### 2.2.2 Index and Search for the Pictures

Well, after finish this part, I have to say it is not as easy as it seems to be at first. It is not because we have to get the url of the picture, the url of the website it's on or the title of the web page, but is the difficulties to get information, or contents of the pictures. Since the structure of the website is quite complex, I really took some time to get enough information I need to search certain pictures.

I choose "http://www.ommoo.com/" to be the page I'm going to index, which is a website offering pictures of the desktop of your computer. So after analyzing the structure of the website for a really hard time, I get the following codes at last, which can make a quite exact index of these pictures.

By the way, if you want to use the index program on other websites, you will

18

Figure 13: the outcome of the site part(3)

have to re-analyse the structure of that site and modify some of the variables in the program so that it can fit the target website.

```python
import sys , os, lucene , threading , time , chardet , urllib2 , re
from datetime import datetime
from BeautifulSoup import BeautifulSoup
from ctypes import *
import urllib
import Queue
import urlparse

"""
This class is loosely based on the Lucene (java implementation) demo class
org.apache.lucene.demo.IndexFiles .   It will take a directory as an argument
and will index all of the files in that directory and downward recursively .
It will index on the file path , the file name and the file contents .
The
resulting Lucene index will be placed in the current directory and called
'index '.
"""

class Ticker(object):

    def __init__(self):
        self.tick = True

    def run(self):
        while self.tick:
```

```
25                sys.stdout.write('.')
26                sys.stdout.flush()
27                time.sleep(1.0)
28
29  class IndexFiles(object):
30      """Usage: python IndexFiles <doc_directory>"""
31
32      def __init__(self, root, storeDir, analyzer):
33
34          if not os.path.exists(storeDir):
35              os.mkdir(storeDir)
36          store = lucene.SimpleFSDirectory(lucene.File(storeDir))
37          writer = lucene.IndexWriter(store, analyzer, True,
38                                      lucene.IndexWriter.MaxFieldLength.LIMITED)
39          writer.setMaxFieldLength(1048576)
40          self.indexDocs(root, writer)
41          ticker = Ticker()
42          print 'optimizing index',
43          threading.Thread(target=ticker.run).start()
44          writer.optimize()
45          writer.close()
46          ticker.tick = False
47          print 'done'
48
49      def indexDocs(self, root, writer):
50          for root, dirnames, filenames in os.walk(root):
51              for filename in filenames:
52                  if filename.endswith('.txt'):
53                      continue
54                  print "adding", filename
55                  try:
56                      path = os.path.join(root, filename)
57                      file = open(path)
58                      buf = file.read()
59                      contents=buf
60                      result = chardet.detect(buf)['encoding']
61                      if result=='GB2312':
62                          contents = buf.decode('gbk').encode('utf8')
63                      file.close()
64                      soup=BeautifulSoup(contents)
65                      url=mydict[filename]
66                      proto, rest = urllib.splittype(url)
67                      site, rest = urllib.splithost(rest)
68                      title=str(soup.head.title.string.strip()).decode('utf8')
69                      flag2=0
70                      for i in soup.findAll('img'):
```

```
71          contents=""
72          flag1=0
73          flag3=0
74          try:
75              contents=contents+' '+i['alt']
76          except:
77              pass
78          tempurl=i['src']
79          imgurl=urlparse.urljoin(url,tempurl)
80          temp=i.parent.parent
81          try:
82              photoid=temp.find('a')['data-photo-id']
83              flag1=1
84          except:
85              pass
86          try:
87              picid=temp.parent.find('article')['id']
88              flag3=1
89          except:
90              pass
91          try:
92              for t in temp.findAll('b'):
93                  try:
94                      contents=contents+' '+t.string.strip()
95                  except:
96                      pass
97          except:
98              pass
99          try:
100             for k in temp.findAll('p'):
101                 try:
102                     contents=contents+' '+k.string.strip()
103                 except:
104                     pass
105         except:
106             pass
107         try:
108             for j in temp.findAll('span',{'class':'title'}):
109                 try:
110                     contents=contents+' '+j.string.strip()
111                 except:
112                     pass
113         except:
114             pass
115         if flag1==1:
116             timetowait=0
```

```python
117                                    try:
118                                        for p in temp.parent.findAll('div',{'class':'car
119                                            if timetowait<flag2:
120                                                timetowait+=1
121                                                continue
122                                            contents=contents+' '+p.string.strip()
123                                            flag2+=1
124                                            break
125                                    except:
126                                        pass
127                                if flag3==1:
128                                    try:
129                                        for q in temp.parent.findAll('div',{'class':'pos
130                                            r=q.find('h1')
131                                            contents=contents+' '+str(r.string).decode('
132                                            break
133                                    except:
134                                        pass
135                                contents=contents.strip()
136                                doc = lucene.Document()
137                                doc.add(lucene.Field("imgurl", imgurl,
138                                                     lucene.Field.Store.YES,
139                                                     lucene.Field.Index.NOT_ANALYZED))
140                                doc.add(lucene.Field("url", url,
141                                                     lucene.Field.Store.YES,
142                                                     lucene.Field.Index.NOT_ANALYZED))
143                                doc.add(lucene.Field("title", title,
144                                                     lucene.Field.Store.YES,
145                                                     lucene.Field.Index.NOT_ANALYZED))
146                                if len(contents) > 0:
147                                    dll=cdll.LoadLibrary("F:\\ICTCLAS50_Windows_32_C\ICT
148                                    dll.ICTCLAS_Init(c_char_p("F:\\ICTCLAS50_Windows_32_C
149                                    strlen = len(c_char_p(contents).value)
150                                    t =c_buffer(strlen*6)
151                                    bSuccess = dll.ICTCLAS_ParagraphProcess(c_char_p(con
152                                    contents=t.value.decode('gbk').encode('utf8')
153                                    ##list=t.value.split()
154                                    ##print ' '.join(list)
155                                    dll.ICTCLAS_Exit()
156                                    doc.add(lucene.Field("contents", contents,
157                                                         lucene.Field.Store.NO,
158                                                         lucene.Field.Index.ANALYZED))
159                                else:
160                                    print "warning: no content in part of %s" % filename
161                                writer.addDocument(doc)
162                        except Exception, e:
```

```
163                              print "Failed_in_indexDocs:", e
164
165  if __name__ == '__main__':
166  ##      if len(sys.argv) < 2:
167  ##          print IndexFiles.__doc__
168  ##          sys.exit(1)
169      lucene.initVM()
170      print 'lucene', lucene.VERSION
171      start = datetime.now()
172      dic= open('F:\\html\index.txt')
173      d = dic.readlines()
174      dic.close()
175      mydict = {}
176      for word in d:
177          key = word.split(';')[0]
178          value = word.split(';')[1]
179          mydict[key] = value
180      try:
181  ##          IndexFiles(sys.argv[1], "index", lucene.WhitespaceAnalyzer(lucene.Vers
182          IndexFiles('F:\\html', "F:\\imgindex", lucene.WhitespaceAnalyzer(lucene.
183          end = datetime.now()
184          print end - start
185      except Exception, e:
186          print "Failed:_", e
```

And after finishing the index part, it's pretty easy to accomplish the rest part.

The code of the searching part is as follows, as it's easy, I won't explain it explicitly.

```
1   from lucene import \
2       QueryParser, IndexSearcher, WhitespaceAnalyzer, SimpleFSDirectory, File, \
3       VERSION, initVM, Version
4
5
6   """
7   This script is loosely based on the Lucene (java implementation) demo class
8   org.apache.lucene.demo.SearchFiles.  It will prompt for a search query, then it
9   will search the Lucene index in the current directory called 'index' for the
10  search query entered against the 'contents' field.  It will then display the
11  'path' and 'name' fields for each of the hits it finds in the index. Note that
12  search.close() is currently commented out because it causes a stack overflow in
13  some cases.
14  """
15  def run(searcher, analyzer):
16      while True:
```

```
17              print
18              print "Hit␣enter␣with␣no␣input␣to␣quit."
19              command = raw_input("Query:")
20              if command == '':
21                  return
22              print
23              print "Searching␣for:", command
24              query = QueryParser(Version.LUCENE_CURRENT, "contents",
25                                      analyzer).parse(command)
26              scoreDocs = searcher.search(query, 50).scoreDocs
27              print "%s␣total␣matching␣documents." % len(scoreDocs)
28
29              for scoreDoc in scoreDocs:
30                  doc = searcher.doc(scoreDoc.doc)
31                  print 'title:', doc.get("title"), 'url:',doc.get("url"), 'imgurl:',
32
33
34  if __name__ == '__main__':
35      STORE_DIR = "F:\\imgindex"
36      initVM()
37      print 'lucene', VERSION
38      directory = SimpleFSDirectory(File(STORE_DIR))
39      searcher = IndexSearcher(directory, True)
40      analyzer = WhitespaceAnalyzer(Version.LUCENE_CURRENT)
41      run(searcher, analyzer)
42      searcher.close()
```

And here are also some screenshots of the files I crawled in the target file 'F:\imgindex', and the pictures of the searching outcomes.

# 3 The Problems I Met in the Experiment and My Solution

Well, there are so many problems I have met in the experiment.
In the first part, I find that sometimes I can't get the content of the pages for the node problem, so I just use 'chardet' in python to find what the code method is and then transform them all to utf8. And then I find a problem in establish the "F:\html" file, and then I find these codes and add them to my program:

```
1  if not os.path.exists("F:\\html"):
2      os.mkdir("F:\\html")
```

These codes mean that if "F:\html" doesn't exist, then it will be set up. And then comes the problem of saving the url and filenames of each page. To make it easy, I save each pair of them in the same line us ";" to separate them from
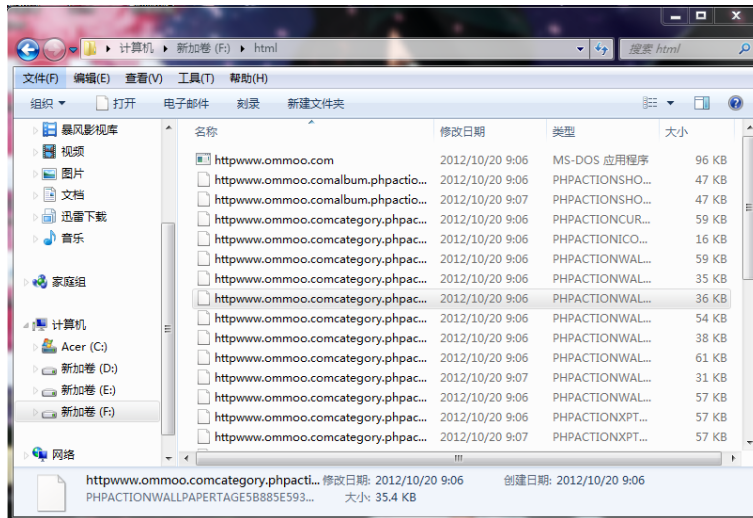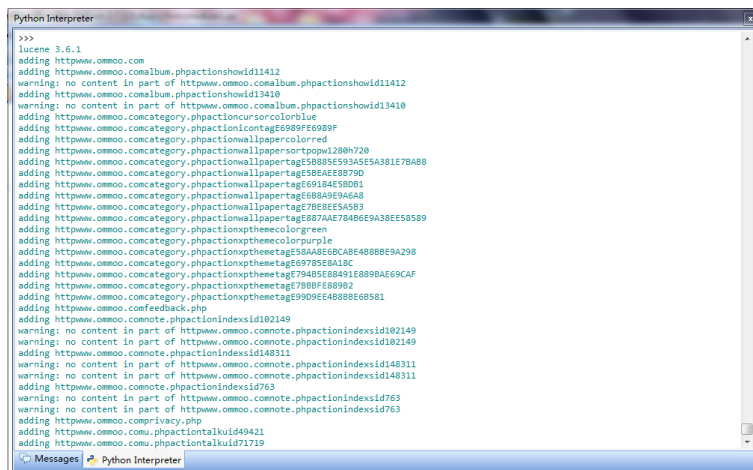
24

Figure 14: the files we get in 'F:\html'



Figure 15: the outcomes of the indexing part

25

Figure 16: the outcomes of the searching part(1)



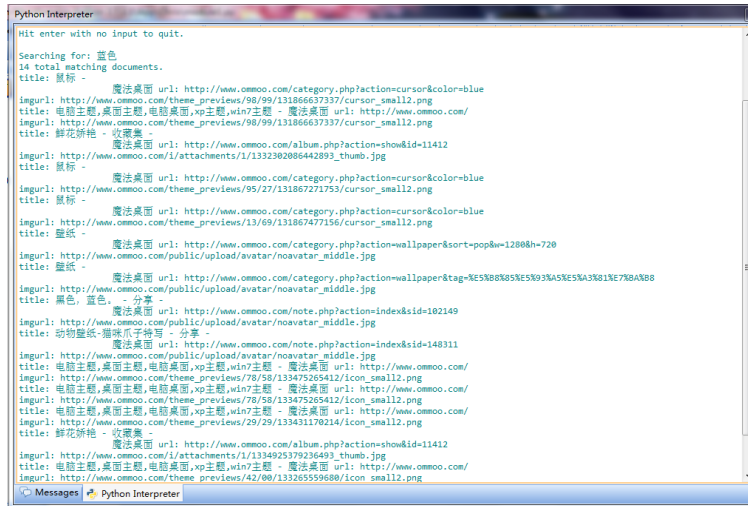Figure 17: the outcomes of the searching part(2)

Figure 18: the outcomes of the searching part(3)

each other. And in this experiment, I choose "www.hao123.com" as the seed web site since I think it can link to plenty of web pages.

As for the second part, I think it is the most complex part in my experiment. I this part, I have to manage to save the name, path, url, title and content of each document in F:\html, except index.txt. And since the method of getting path ,name and the content of the document has been given in the demo, so I only need to get the url and title. And in order to get the url quickly and correctly, I will use the 'index.txt' I set up in step 1. And in this part, I use the following codes to change it into a dictionary.

```
1  dic= open('F:\\html\index.txt')
2  d = dic.readlines()
3  dic.close()
4  mydict = {}
5  for word in d:
6  key = word.split(';')[0]
7  value = word.split(';')[1]
8  mydict[key] = value
```

And in this way we make a dictionary 'mydict', making it easy to find the corresponding url of each page. Then I use BeautifulSoup to find the title using this code:

```
1  title=str(soup.head.title.string).decode('utf8')
```

It seems that this part is finished, but actually, it is far from saying so now. If we read the third part carefully, we will find that we have to separate the

27

Chinese words using some dictionaries ad then pass them to the analyzer in lucene. Without this ,we will find the outcome of our search will be a mass. So I decide to use ICTCLAS50 to separate the Chinese words with these codes:

```
1  from ctypes import *
2  dll=cdll.LoadLibrary("F:\\ICTCLAS50_Windows_32_C\ICTCLAS50.dll")
3  dll.ICTCLAS_Init(c_char_p("F:\\ICTCLAS50_Windows_32_C"))
4  strlen = len(c_char_p(contents).value)
5  t =c_buffer(strlen*6)
6  bSuccess = dll.ICTCLAS_ParagraphProcess(c_char_p(contents),c_int(strlen),t,c_int
7  contents=t.value.decode('gbk').encode('utf8')
8  ##list=t.value.split()
9  ##print ' '.join(list)
10 dll.ICTCLAS_Exit()
```

When I finish the separating work, I find I couldn't find any result in the third part. Then I realized that it is because I didn't change the contents splited by the dictionary to utf8 code, since ICTCLAS50 can support gbk code, which means maybe my contents are just in gbk form. And then when I changed the code, the result is ok. And I need to say, in order to make the content easily to be checked later, I use SimpleAnalyzer to separate in the following part.

In the third part, it is quite easy since the most difficult parts have been solved earlier. And what I need to do is to change the StandardAnalyzer in the demo to SimpleAnalyzer, and make sure the outcome contains path, title, url and name.

Here is a screenshot of where I put those files.

In the second half, I happened to find a method that can get the domain name of the website quickly on the Internet. It is attained by using the urllib library. And the codes are as follows:

```
1  soup=BeautifulSoup(contents)
2  url=mydict[filename]
3  proto, rest = urllib.splittype(url)
4  site, rest = urllib.splithost(rest)
```

And in this way we can get 'site' as the domain name.

And as for the image search, I decided to use 'www.taobao.com' at first, but later I found there is some problem with getting the url of the pictures. In this condition, some pictures have the url form as 'data-ks-lazyload' instead of the normal 'src' form. I looked it up and found it was because the so-called 'delay loading of the picture'. And as time presses, I hardly have any time to study it deeply to solve the problem, so I have to change the target website. Maybe I will try to solve it later.
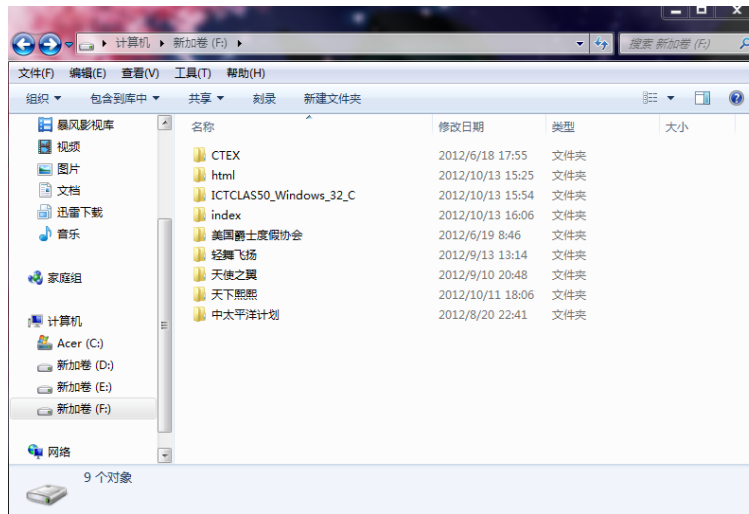
Figure 19: the files we need(in English)

# 4  Some of My Thoughts

From lab1 to lab3, it's not hard to find that we are nearer to the search engine step by step. We learn how to get urls in lab1, how to crawl the pages in lab2 and how to search according to the query after making an index in this experiment. So I really think it's fun to study in this field. Though I have some tough time learning things I never heard of and spend tons of time programming. But when I see the achievement from my own hands, I feel real happiness that can't be expressed with words.

And as for the second half, we know that it can be more exactly when we try to search some information on a certain website. And I do think it is of great use because most of the times, we really want to search something on one page, and our experiment enables us to achieve it.

In the experiment, I have ever forgot to save the successful code and lost them, taking me some time to write them again. But luckily, in this process, I think I am better at commanding the principle of the searching method. But I think the most annoying thing is that python is really slow when running. If it can be a little faster, I think it will be a more relaxed thing searching information using the program written of our own.