# Report of Lab4

Yuxiang Chen 5110309783

November 4, 2012

## Contents

## 1   The Purpose of Lab4 and My Preparation

It is definitely sure that we have to build our own search engine by ourselves. But how to write such a program that enables us to search certain information in a graphic user interface, e.g. in a web page(port)? The framework 'web.py' gives us the chance.

In this experiment, we need to learn how to use this framework to visit a certain server where we can do our searching jobs. To accomplish this experiment, I look up many things in the website introducing web in python, and find that 'web.py' is just a framework that makes it possible to connect our searching programs with the html lists. And there are also some other things to do, such as print out the url, main contents and the title of the web page in the form of a hyperlink.

## 2   The Main Part of the Experiment

To make it easy to understand, I split this experiment into three parts, which consists of setting up an index, searching for information you required and establishing the framework of the html page to be presented later.

## 2.1 Making an Index First

This part is mainly the same as the experiments we did earlier, except that we need to do some work to analyse and store the text in the neighborhood of the keyword requested. And we also need to exclude the useless messages in the content of the page, such as how the web site set up the html page. And I will present them in the following codes:

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

def filter_tags(htmlstr):
    re_cdata=re.compile('//<!\[CDATA\[[^>]*//\]\]>',re.I)
    re_script=re.compile('<\s*script[^>]*>[^<]*<\s*/\s*script\s*>',re.I)
    re_style=re.compile('<\s*style[^>]*>[^<]*<\s*/\s*style\s*>',re.I)
    re_br=re.compile('<br\s*?/?>')
    re_h=re.compile('</?\w+[^>]*>')
    re_comment=re.compile('<!--[^>]*-->')
    s=re_cdata.sub('',htmlstr)
    s=re_script.sub('',s)
    s=re_style.sub('',s)
    s=re_br.sub('\n',s)
    s=re_h.sub('',s)
    s=re_comment.sub('',s)
    blank_line=re.compile('\n+')
    s=blank_line.sub('\n',s)
    s=replaceCharEntity(s)
    return s


def replaceCharEntity(htmlstr):
    CHAR_ENTITIES={'nbsp':' ','160':' ','lt':'<','60':'<','gt':'>','62':'>','amp
    re_charEntity=re.compile(r'&#?(?P<name>\w+);')
    sz=re_charEntity.search(htmlstr)
    while sz:
        entity=sz.group()
        key=sz.group('name')
        try:
            htmlstr=re_charEntity.sub(CHAR_ENTITIES[key],htmlstr,1)
            sz=re_charEntity.search(htmlstr)
        except KeyError:
            htmlstr=re_charEntity.sub('',htmlstr,1)
```

```
39              sz=re_charEntity.search(htmlstr)
40      return htmlstr
41
42
43  def repalce(s,re_exp,repl_string):
44      return re_exp.sub(repl_string,s)
45
46
47  class Ticker(object):
48
49      def __init__(self):
50          self.tick = True
51
52      def run(self):
53          while self.tick:
54              sys.stdout.write('.')
55              sys.stdout.flush()
56              time.sleep(1.0)
57
58  class IndexFiles(object):
59      """Usage: python IndexFiles <doc_directory>"""
60
61      def __init__(self, root, storeDir, analyzer):
62
63          if not os.path.exists(storeDir):
64              os.mkdir(storeDir)
65          store = lucene.SimpleFSDirectory(lucene.File(storeDir))
66          writer = lucene.IndexWriter(store, analyzer, True,
67                                      lucene.IndexWriter.MaxFieldLength.LIMITED)
68          writer.setMaxFieldLength(1048576)
69          self.indexDocs(root, writer)
70          ticker = Ticker()
71          print 'optimizing index',
72          threading.Thread(target=ticker.run).start()
73          writer.optimize()
74          writer.close()
75          ticker.tick = False
76          print 'done'
77
78      def indexDocs(self, root, writer):
79          for root, dirnames, filenames in os.walk(root):
80              for filename in filenames:
81                  if filename.endswith('.txt'):
82                      continue
83                  print "adding", filename
84                  try:
```

3

```python
                        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')
                        new_contents=filter_tags(contents)
                        new_contents=str(new_contents).strip().decode('utf8')
                        pos=new_contents.find('>')
                        new_contents=new_contents[pos+1:]
                        temp=new_contents.split()
                        newtext=' '.join(temp)
                        contents=''.join(soup.findAll(text=True))
                        doc = lucene.Document()
                        doc.add(lucene.Field("text", newtext,
                                            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(content
                            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,
                                                lucene.Field.Index.ANALYZED))
                        else:
                            print "warning: no content in %s" % filename
                        writer.addDocument(doc)
                except Exception, e:
                    print "Failed in indexDocs:", e
```

```
131
132  if __name__ == '__main__':
133  ##      if len(sys.argv) < 2:
134  ##          print IndexFiles.__doc__
135  ##          sys.exit(1)
136      lucene.initVM()
137      print 'lucene', lucene.VERSION
138      start = datetime.now()
139      dic= open('F:\\html\index.txt')
140      d = dic.readlines()
141      dic.close()
142      mydict = {}
143      for word in d:
144          value=''
145          try:
146              key = word.split(';')[0]
147              value = word.split(';')[1]
148              mydict[key] = value
149          except:
150              pass
151      try:
152  ##          IndexFiles(sys.argv[1], "index", lucene.SimpleAnalyzer(lucene.Version..
153              IndexFiles('F:\\html', "F:\\index", lucene.SimpleAnalyzer(lucene.Version
154              end = datetime.now()
155              print end - start
156      except Exception, e:
157              print "Failed: ", e
```

So after this part, most of the useless information in the web page is eliminated, leaving the messages we need in the 'text' variable. There's any output after the program being processed. So only the screenshot of the index files will be given.

## 2.2   The Main Searching Part of the Program

In this part, I will show you the main process of searching according to a user's query. And I try my best to make it as flexible as it could. E.g., if you search for a certain Chinese word, the program will search it by treating the word as a whole expression and then treating it as maybe two words by split them. So in this way, I try to return as much information as I can. And for other potential errors, I deal with the 'try...except' sentence.

```
1  import web
2  from web import form
3  import urllib2
```
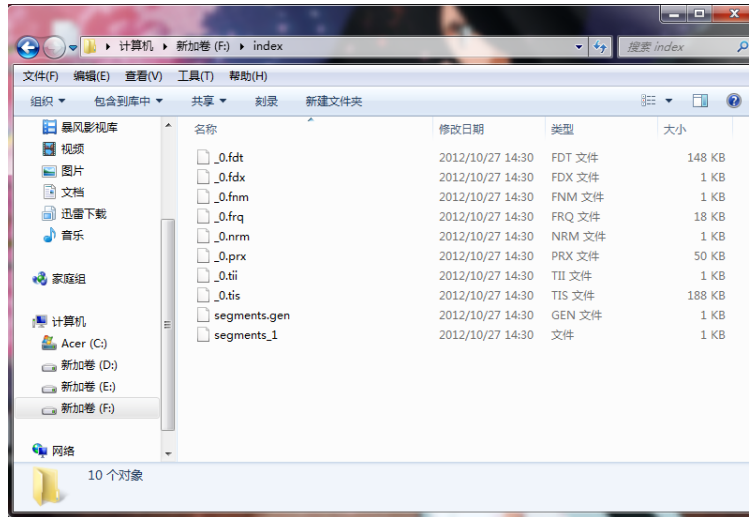
Figure 1: the files of the index

```
 4  import os
 5  from lucene import \
 6      QueryParser, IndexSearcher, SimpleAnalyzer, SimpleFSDirectory, File, \
 7      VERSION, initVM, Version
 8  from ctypes import *
 9
10  urls = (
11      '/', 'index',
12      '/s', 's'
13  )
14
15
16  render = web.template.render('templates') # your templates
17
18  login = form.Form(
19      form.Textbox('keyword'),
20      form.Button('Search'),
21  )
22
23  def func(command, searcher, analyzer):
24      if len(command) > 0:
25          dll=cdll.LoadLibrary("F:\\ICTCLAS50_Windows_32_C\ICTCLAS50.dll")
26          dll.ICTCLAS_Init(c_char_p("F:\\ICTCLAS50_Windows_32_C"))
27          strlen = len(c_char_p(command).value)
28          t =c_buffer(strlen*6)
29          bSuccess = dll.ICTCLAS_ParagraphProcess(c_char_p(command),c_int(strlen),
```

6

```
30          command=t.value.decode('gbk').encode('utf8')
31          ##list=t.value.split()
32          ##print ' '.join(list)
33          dll.ICTCLAS_Exit()
34          command=command.decode('utf8')
35      query = QueryParser(Version.LUCENE_CURRENT, "contents",analyzer).parse(comma
36      scoreDocs = searcher.search(query, 50).scoreDocs
37      total=len(scoreDocs)
38      qtitle=[]
39      qurl=[]
40      qnew_text1=[]
41      qnew_text2=[]
42      q_query=[]
43      new_query=str(query).replace ('contents:','').decode('utf8')
44      if total==0:
45          return command,qurl,qtitle,qnew_text1,qnew_text2,q_query,total
46      for scoreDoc in scoreDocs:
47          doc = searcher.doc(scoreDoc.doc)
48          text=doc.get("text")
49          new_text=str(text).decode('utf8')
50          temp_query=new_query.replace('_','')
51          num=new_text.find(temp_query)
52          query_len=len(temp_query)
53          splited_query=new_query.split('_')
54          splitlen=len(splited_query)
55          if (num!=-1):
56              try:
57                  new_text1=str(new_text[num-30:num]).strip().decode('utf8','ignor
58                  new_text2=str(new_text[num+query_len:num+30+query_len]).strip().
59              except:
60                  try:
61                      new_text1=str(new_text[num-30:num]).strip().decode('utf8','i
62                      new_text2=""
63                  except:
64                      try:
65                          new_text1=""
66                          new_text2=str(new_text[num+query_len:num+30+query_len]).
67                      except:
68                          new_text1=""
69                          new_text2=""
70              q_query.append(temp_query)
71          else:
72              for i in range(splitlen):
73                  parted_query=splited_query[i].decode('utf8')
74                  num=new_text.find(parted_query)
75                  if num==-1:
```

7

```
76                          continue
77                      query_len=len(parted_query)
78                      try:
79                          new_text1=str(new_text[num-30:num]).strip().decode('utf8','i
80                          new_text2=str(new_text[num+query_len:num+30+query_len]).stri
81                      except:
82                          try:
83                              new_text1=str(new_text[num-30:num]).strip().decode('utf8
84                              new_text2=""
85                          except:
86                              try:
87                                  new_text1=""
88                                  new_text2=str(new_text[num+query_len:num+30+query_le
89                              except:
90                                  new_text1=""
91                                  new_text2=""
92                      q_query.append(parted_query)
93                      break
94                  if num==-1:
95                      total=total-1
96                      continue
97              title=doc.get("title")
98              url=doc.get("url")
99              qurl.append(url)
100             qtitle.append(title)
101             qnew_text1.append(new_text1)
102             qnew_text2.append(new_text2)
103         return command,qurl,qtitle,qnew_text1,qnew_text2,q_query,total
104
105 class index:
106     def GET(self):
107         f = login()
108         return render.formtest(f)
109
110 class s:
111     def GET(self):
112         vm_env = initVM()
113         form1 = login()
114         user_data = web.input()
115         vm_env.attachCurrentThread()
116         STORE_DIR = "F:\\index"
117         directory = SimpleFSDirectory(File(STORE_DIR))
118         searcher = IndexSearcher(directory, True)
119         analyzer = SimpleAnalyzer(Version.LUCENE_CURRENT)
120         a,b,c,d,e,f,g= func(user_data.keyword,searcher,analyzer)
121         searcher.close()
```

```
122            return render.result(form1,a,b,c,d,e,f,g)
123
124  if __name__ == "__main__":
125      while True:
126          app = web.application(urls, globals())
127          app.run()
```

Of course, the web pages to be presented afterwards are also included in the
above part, but how to set up the web page will be presented in the next part.
And in order to make the pages a little more beautiful, I revised some of the
variables when necessary. And the other parts are nearly the same as the 'search'
program in lab3.
The pictures will be given after I finished introducing the 'html' part.

## 2.3   The Part of Establishing the Web Page

In this part, I wrote two html files for the searching web page, and as required,
it can present the url, some contents near the keywords, and title in the form
of a hyperlink of the result. Here are the codes of the two pages:
The first is the formtest.html:

```
1  $def with(form)
2  <head>
3  <title> C h r i s          </title>
4  <form name="input" form action="/s" method="GET">          :
5  <input type="keyword" name="keyword" />
6  <input type="submit" value="          " />
7  </form>
```

There are some parts written in Chinese which can't be showed in the pdf file,
you can look for it in my code files. And the following is the result.html:

```
1  $def with (form,name,qurl,qtitle,qnew_text1,qnew_text2,q_query,total)
2
3  <title> C h r i s          </title>
4  <form name="input" form action="/s" method="GET">          :
5  <input type="keyword" name="keyword" />
6  <input type="submit" value="          " />
7  </form>
8  <head><h1><strong></strong><font color="purple">Chris  L u n e          </h1><strong
9              <h2><font color="black">                    "_$name_":</h2>
10             <h4><font color="black">          $total    </h4>
11                  <h4><b><font color="black">
12 $if qtitle and qurl and len(qnew_text1) and len(qnew_text2):
13         $for i in range(total):
14         <h3><a href="_$qurl[i]_"><font color="blue">$qtitle[i]</a></h3>
15         <div><i><font color="black">$qnew_text1[i].strip()<font color="red">$q_q
16         <h5><div><font color="green"><i>$qurl[i]</i></div></h5>
```

9

```
17 │ $else:
18 │     <em>There's no result!</em>!
```

So after these preparations, the result of the searching program is shown in the following screenshots:



Figure 2: the outcome 1



Figure 3: the outcome 2

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

As for the first part, the main problem I meet is how to eliminate the useless information near the key words, and after looking for some solutions on the

Figure 4: the outcome 3



Figure 5: the outcome 4



Figure 6: the outcome 5

Internet, I find the following method to solve it:

```
 1  import re
 2  def filter_tags(htmlstr):
 3      re_cdata=re.compile('//<!\[CDATA\[[^>]*//\]\]>',re.I)
 4      re_script=re.compile('<\s*script[^>]*>[^<]*<\s*/\s*script\s*>',re.I)
 5      re_style=re.compile('<\s*style[^>]*>[^<]*<\s*/\s*style\s*>',re.I)
 6      re_br=re.compile('<br\s*?/?>')
 7      re_h=re.compile('</?\w+[^>]*>')
 8      re_comment=re.compile('<!--[^>]*-->')
 9      s=re_cdata.sub('',htmlstr)
10      s=re_script.sub('',s)
11      s=re_style.sub('',s)
12      s=re_br.sub('\n',s)
13      s=re_h.sub('',s)
14      s=re_comment.sub('',s)
15      blank_line=re.compile('\n+')
16      s=blank_line.sub('\n',s)
17      s=replaceCharEntity(s)
18      return s
19
20
21  def replaceCharEntity(htmlstr):
22      CHAR_ENTITIES={'nbsp':' ','160':' ','lt':'<','60':'<','gt':'>','62':'>','amp
23      re_charEntity=re.compile(r'&#?(?P<name>\w+);')
24      sz=re_charEntity.search(htmlstr)
25      while sz:
26          entity=sz.group()
27          key=sz.group('name')
28          try:
29              htmlstr=re_charEntity.sub(CHAR_ENTITIES[key],htmlstr,1)
30              sz=re_charEntity.search(htmlstr)
31          except KeyError:
32              htmlstr=re_charEntity.sub('',htmlstr,1)
33              sz=re_charEntity.search(htmlstr)
34      return htmlstr
```

Then in this way, we can change some marks into the punctuations they present,
and also delete some other messages such as the tag of the page, etc.
And in the second searching part, there are to many errors to deal with, and I
write all solutions to the conditions I have met in the 'try...except' clause. Also,
in the last part, I firstly have some trouble with giving out all results using a
circling clause, then I learn how to solve it from the reference web site. And
below are the codes:

```
 1  $for i in range(total):
```

```
2         <h3><a href="_$qurl[i]_"><font color="blue">$qtitle[i]</a></h3>
3         <div><i><font color="black">$qnew_text1[i].strip()<font color="red">$q_q
4         <h5><div><font color="green"><i>$qurl[i]</i></div></h5>
```

And after solving these problems ,I basically finish the searching program.

## 4   Some of My Thoughts

In lab4, it's easy to find that we are able to write some part of our own searching engine by ourselves. Although at first I'm a little unfamiliar with how to write some clause in the html file, I manage to solve it by learning from the Internet. And now I feel some kind of satisfied by searching for some queries using my own searching files. And now that we have finished the searching program of texts in the page, I think it's also not hard to search and return some images in the same way.

But there may still exist some mistakes which I haven't found in my program yet, so I'll still testing it in the following days, trying to make it as perfect as I want.