

2014-05-14 项目二：miniSearch项目文档

项目文档

说明

- 开发环境：Linux、c++、PHP、TCP、多进程
- 程序框架：客户端通过PHP页面用TCP协议发送搜索词到服务器，服务器接收后交予子进程处理，子进程进行切词，计算权重并与网页库进行相似度计算，逆序输出相似度最高的N篇文章。

第一阶段：搭建框架 0514

1.目录结构

1.	bin	//存放可执行文件和执行脚本
2.	conf	//配置文件，包括服务器地址，网页库、网页库索引、单词索引、停用词表的路径等
3.	data	//存放网页库、网页库索引、单词索引、停用词表
4.	include	//存放头文件
5.	log	//存放系统日志文件
6.	src	//存放 cpp 代码文件
7.	Makefile	//Makefile

2.将服务器设置为守护进程

1.	static void Daemon() {
2.	const int MAXFD = 64;
3.	if (fork() != 0) //父进程退出
4.	exit(0);
5.	setsid(); //成为新进程组组长和新会话领导，脱离控制终端
6.	//chdir("/"); //此时目录不能改变，因为代码内文件路径是相对路径
7.	umask(0); //重设文件访问权限掩码
8.	for (int i = 0; i < MAXFD; i++) //尽可能关闭所有从父进程继承来的文件
9.	{
10.	if (i == 1) //关闭所有文件描述符，除了1，以便通过脚本运行时将输出流重定向至./log/log.txt
11.	{
12.	continue;
13.	}
14.	close(i);
15.	}
16.	}

3.封装tcp类

4.搭建多进程框架

每当检测到客户端连接时，利用 fork() 生成一个子进程，让子进程去处理查询任务，并将查询结果发回客户端。

这里涉及到子进程的回收问题。若利用父进程的 waitpid() 回收，有可能造成父进程阻塞，不能及时处理新的访问，或者有可能产生僵尸进程，造成资源浪费。解决方法是：在 fork() 函数前设置 signal(SIGCHLD, wait_child); 或 signal(SIGCHLD, SIG_IGN);，让系统接管子进程的回收。

1.	static void wait_child(int pid) {
2.	if ((pid = waitpid(-1, NULL, WNOHANG)) > 0) {
3.	cout << "pid: " << pid << " exit!" << endl;
4.	}
5.	}
6.	int main(int argc, char **argv) {
7.	
8.	Daemon(); //守护进程
9.	
10.	std::ifstream ifs("./conf/config.txt");
11.	std::string line, conf, ip, port;
12.	while (getline(ifs, line)) {
13.	istringstream iss(line);
14.	iss >> conf;
15.	if (conf == "address:") {
16.	iss >> ip >> port;
17.	break;
18.	}
19.	}
20.	#ifndef DEBUG //输出ip和端口
21.	std::cout << ip << " " << port << std::endl;
22.	#endif
23.	
24.	signal(SIGCHLD, wait_child); //通过wait_child(int)函数来等待回收子进程退出后回收子进程的资源(也可使用下面这句语句回收子进程资源)
25.	
26.	/* 或单独使用下面的函数 */
27.	//signal(SIGCHLD, SIG_IGN); //父进程忽略子进程退出时传回的信号，让子进程结束的信号由系统接收，由1号进程负责回收子进程的资源
28.	
29.	Task task(50,10); //Task(最大显示的篇数，每篇显示的行数)
30.	TcpSocket mysocket(ip, port);
31.	#ifndef DEBUG //打印标志信息
32.	cout << "=== ready to recv message ===" << endl;
33.	#endif

```
34.     while (int client_fd = mysocket.accept_connection()) {
35.         //当检测到新的连接时，创建子进程来处理任务
36.         struct sockaddr_in client_addr = mysocket.get_client_addr();
37.         int pid = fork();
38.         if (pid == 0) {
39.             while (true) {
40.
41.                 //recv
42.                 std::string recv_buf;
43.                 mysocket.recv_message(client_fd, recv_buf);
44.                 cout << "recv:" << recv_buf << endl;
45.                 if (recv_buf == "exit") {
46.                     cout << "client " << "ip: "
47.                         << inet_ntoa(client_addr.sin_addr) << " "
48.                         << ntohs(client_addr.sin_port) << " exit" << endl;
49.                     close(client_fd);
50.                     exit(0);
51.                 }
52.
53.                 //send
54.                 std::string send_buf;
55.                 send_buf = task.search(std::string(recv_buf));
56.                 mysocket.send_message(client_fd, send_buf);
57.             }
58.         } else {
59.             //close(client_fd);
60.             //waitpid(pid, NULL, WNOHANG); //signal()之后不需要wait()
61.         }
62.     }
63.     return 0;
64. }
```

第二阶段：网页库构建 0515

遍历目录读库文件，拼接成标准格式，然后写入文件，并同时建立库索引

```
1.  /*
2.  * 网页库文件格式:
3.  * <doc><docid>1</docid>
4.  * <url>http://baidu.com/</url>
5.  * <title>标题</title>
6.  * <content>标题 + 内容</content></doc>
7.  */
8. doc = "<doc><docid>" + string(id) + "</docid><url>" +
9.     string(entry->d_name) + "</url>" + "<title>" + title +
10.     "</title><content>" + content + "</content></doc>\n";
11. lib_vec.pushback(doc); //每篇doc存入vector
```

```
1.  /*
2.  * 写入网页库文件的同时，
3.  * 记录每篇doc的文件偏移量，和该篇doc的长度，
4.  * 写入库索引文件
5.  */
6. for(vector<string>::iterator iter = lib_vec.begin(); iter != lib_vec.end(); ++iter)
7. {
8.     ofs_index << i << " " << ofs_lib.tellp();
9.     //向index写入索引 docid, start_pos
10.    ofs_lib << *iter;
11.    //向lib写入doc内容
12.    ofs_index << " " << (*iter).size() << endl;
13.    //向index写入索引 size
14.    i++;
15. }
```

编码的转换

由于语料库是gbk编码，在Ubuntu Terminal显示是乱码，并且由于网页输入是utf8编码，而编码转换效率较低而且容易出现段错误，在构建网页库之前，为便于测试，将所有语料库批量装换为utf8编码。转换时，会出现类似"^M"的乱码，这是由于windows下的回车是"\r\n"，在Ubuntu下显示的问题。必须去除该字符，因为者会造成某些语句被下一行语句覆盖！！！

因为\r是回到行首，而\n才是换行。

第三阶段：网页去重 0516

遍历网页库文件，分别取出每一篇doc，切词，去停用词后统计每篇doc的词频，放入优先级队列。从优先级队列中取出k个词频最高的单词，存入

```
1.     map<int, map<string, int> > doc_feature;
2.     //docid, word, freq
```

它可以代表每篇doc的特征。

然后开始去重：

```
1.     //去重
2.     int * arr = new int [doc_feature.size()]; //构造一个数组
3.     for(int i = 0; i != (int)doc_feature.size(); ++i)
4.     {
5.         arr[i] = i+1; //给数组赋初始值
6.     }
7.     int ix = 0, iy = 0;
```

```
8. map<int, map<string, int> >::iterator iter_end = doc_feature.end();
9. iter_end--;
10. for(map<int, map<string, int> >::iterator iter_x = doc_feature.begin(); iter_x != iter_end; ++iter_x) //待比较文章x
11. {
12.     if(arr[ix] == 0) //当前文章x已经被去除，跳过
13.     {
14.         iter_x ++;
15.         ix ++;
16.         continue;
17.     }
18.     map<int, map<string, int> >::iterator iter_y = iter_x;
19.     iter_y++;
20.     iy = ix; iy++;
21.     for(; iter_y != doc_feature.end(); ++iter_y) //待比较文章y
22.     {
23.         if(arr[iy] == 0) //当前文章y已经被去除，跳过
24.         {
25.             iter_y ++;
26.             iy ++;
27.             continue;
28.         }
29.         if(compare_two_doc(iter_x->second, iter_y->second) >= 6) //重复
30.         {
31.             #ifndef NDEBUG //有重复时输出显示发生重复的两个docid和去除的docid
32.                 cout << "repeat! " << iter_x->first << " " << iter_y->first << "   trim " << arr[iy] << "   total(" << doc_feature.size() << endl;
33.             #endif
34.             arr[iy] = 0; //直接去除docid大的一篇
35.         }
36.         iy++;
37.     }
38.     ix++;
39. }

//将去重数组写入文件
42. string outfile = argv[1] + string("_dup_arr.dat");
43. std::ofstream dup_arr_file(outfile.c_str());
44. for(int i = 0; i != (int)doc_feature.size(); ++i)
45. {
46.     dup_arr_file << arr[i] << " ";
47. }
48. dup_arr_file.close();
49. delete [] arr;
```

比较判断两篇 doc 是否重复的 top k 策略：

```
1.  /*
2.   * 比较两篇doc的词频最高的10个单词
3.   * 返回相似度
4.   * 即相同单词的个数
5.   *
6.   * 其他方法：
7.   *   将一个map插入另一个map
8.   *   每插入一个元素时
9.   *       如果相同则插入后size不变
10.  *       如果不同则size +1
11.  */
12. int compare_two_doc(const map<string, int> &mp1, const map<string, int> &mp2)
13. {
14.     int count_dup = 0;
15.     for(map<string, int>::const_iterator iter1 = mp1.begin(); iter1 != mp1.end(); ++iter1)
16.     {
17.         for(map<string, int>::const_iterator iter2 = mp2.begin(); iter2 != mp2.end(); ++iter2)
18.         {
19.             if(iter1->first == iter2->first)
20.                 count_dup++;
21.         }
22.     }
23.     return count_dup;
24. }
```

根据以上去重方法得到的去重数组、和原来的网页库、网页库索引，建立新的网页库和网页库索引。重做新网页库的docid。

第四阶段：建索引 0517

目标：

```
1. //单词 docid 归一化之后的权重
2. word1 docid normalized_power docid normalized_power
3. word2 docid normalized_power ...
4. ...
```

第一步，根据公式求权重：

$$power_{word} = tf_{doc} * log(\frac{N}{df_{word}})$$

```
1. power = tf * log(N/df);
2. //power 该词在某一篇doc中的权重
3. //tf 该词在某一篇doc中的词频
4. //df 该词的文档频率
5. //N 文档总数
```



```
1.  /*
2.   * 计算权重
3.   */
4. void compute_power(map<string, map<int, int> > &word_docid_freq, map<string, map<int, double> > &word_power)
5. {
6.     cout << "!!!!!!compute_power!!!!!!" << endl;
7.     for(map<string, map<int, int> >::iterator iter = word_docid_freq.begin(); iter != word_docid_freq.end(); ++iter)
8.         //遍历所有单词
9.         {
10.            string word = iter->first;
11. #ifndef DEBUG    //测试输出单词的权值
12.             //cout << word << ": " << endl;
13. #endif
14.             for(map<int, int>::iterator it = word_docid_freq[word].begin(); it != word_docid_freq[word].end(); ++it)
15.                 //遍历含有该单词的docid
16.                 {
17.                     int docid = it->first;
18.                     int freq = it->second;
19.                     word_power[word][docid] = freq*log(0.05 + (double)word_docid_freq.size()/((double)word_docid_freq[word].size()));
20.
21. #ifndef DEBUG    //测试输出单词的权值
22.                     cout << "docid: " << docid << " power: " << word_power[word][docid] << endl;
23. #endif
24.                 }
25.             }
26.             cout << "compute over!!!" << endl;
27.         }
```

第二步，根据公式，把权重归一化：

$$normalized_power_{w1} = \frac{power_{w1}}{\sqrt{\sum_{i=1}^n (power_{wi})^2}}$$

```
1.  //例如对于文档 d1,可以得到一个所有单词的权重序列
2.  w1, w2, w3 ..... wn
3.  //然后需要归一化,即
4.  w1 = w1 / sqrt(pow(w1) + pow(w2) + ..... pow(wn))
```

```
1.  /*
2.   * 权值归一化
3.   */
4. void get_normalized_power(map<string, map<int, double> > &word_power, map<int, map<string, int> > &doc_feature, map<string, map<int, double> > &normalized_power)
5. {
6.     cout << "!!!!compute normalized_power!!!!" << endl;
7.     map<int, double> power_2;    //docid权重的平方和开根号
8.     cout << "!!!compute power^2!!!" << endl;
9.     for(map<int, map<string, int> >::iterator iter_doc = doc_feature.begin(); iter_doc != doc_feature.end(); ++iter_doc)
10.        //遍历docid
11.        {
12.            int docid = iter_doc->first;
13.            cout << "docid: " << docid << endl;
14.            double base = 0;
15.            for(map<string, int>::iterator iter_word = iter_doc->second.begin(); iter_word != iter_doc->second.end(); ++iter_word)
16.                //计算分母 sqrt(pow(w1) + pow(w2) + ... + pow(wn))
17.                {
18.                    double power = word_power[iter_word->first][docid];
19.                    base += (power*power);
20.                }
21.            power_2[docid] = sqrt(base);
22.        }
23.     cout << "compute sqrt(power^2) over" << endl;
24.
25.     cout << "!!!!compute normalized_power!!!!" << endl;
26.     for(map<int, map<string, int> >::iterator iter_doc = doc_feature.begin(); iter_doc != doc_feature.end(); ++iter_doc)
27.        //遍历docid
28.        {
29.            int docid = iter_doc->first;
30.            for(map<string, int>::iterator iter_word = iter_doc->second.begin(); iter_word != iter_doc->second.end(); ++iter_word)
31.            {
32.                normalized_power[iter_word->first][docid] = word_power[iter_word->first][docid]/power_2[docid];
33.                //取出分子，与分母相除，计算归一化权重
34.            }
35.        }
36.     cout << "compute over!!!" << endl;
37. }
```

第三步，单词索引写入文件：

```
1.  ofstream ofs_word_index("words.index");
2.     //写单词索引文件
3.  for(map<string, map<int, double> >::iterator iter = normalized_power.begin(); iter != normalized_power.end(); ++iter)
4.  {
5.      ofs_word_index << iter->first << " ";
6.      for(map<int, double>::iterator it = iter->second.begin(); it != iter->second.end(); ++it)
7.      {
8.          ofs_word_index << it->first << " " << it->second << " ";
9.      }
10.     ofs_word_index << endl;
11. }
```

第五阶段：计算文本相似度 0518

根据公式，计算两篇doc之间的相似度：

$$sim(doc_1, doc_2) = \sum_{i \in doc_1 \cap doc_2} (n_power_doc1_{wi} * n_power_doc2_{wi})$$

```
1. //计算 doc_x, doc_y 之间的相似度
2. //两篇doc之间的公共单词是 word_1, word_2 ... word_n
3. sim(doc_x, doc_y) = (wx_1*wy_1 + wx_2*wy_2 + ... + wx_n*wy_n);
```

第六阶段：查询模块 0518

对输入的查询query进行切词（去停用词），计算出同时包含所有查询词的所有docid：

```
1. /*
2.  * 求出同时包含不同搜索词的docid,
3.  * 放入 set<int> common_docid 中
4.  */
5. static void get_common_docid(map<string, map<int, double> > &m_word_index, vector<string> &words, set<int> &common_docid)
6. {
7.     if(words.size() >= 1)
8.     {
9.         for(map<int, double>::iterator iter = m_word_index[words[0]].begin(); iter != m_word_index[words[0]].end(); ++iter)
10.        {
11.            common_docid.insert(iter->first);
12.            //将包含第一个单词的docid全部输入到set<int>中
13.        }
14.    }
15.    if(words.size() > 1)
16.    {
17.        for(vector<string>::size_type ix = 1; ix != words.size(); ++ix)
18.            //从第二个单词开始遍历, 查找同时含有搜索关键词的docid
19.        {
20.            for(set<int>::iterator iter = common_docid.begin(); iter != common_docid.end(); )
21.                //遍历set<int>查找该docid是否存包含当前搜索词
22.                //这里涉及到遍历set并用erase删除元素
23.                //***需要判断边界!!!! (重要) ***//
24.            {
25.                set<int>::iterator it_back = iter; //备份迭代器
26.                bool is_begin = false;
27.                if(it_back == common_docid.begin())
28.                {
29.                    is_begin = true;
30.                }
31.                else
32.                {
33.                    it_back --; //备份迭代器
34.                }
35.
36.                if(!m_word_index[words[ix]].count(*iter))
37.                    //set<int>的docid不包含当前搜索词
38.                {
39.                    //cout << "not common " << *iter << endl;
40.                    common_docid.erase(iter);
41.                    //删除元素 (docid)
42.                    if(is_begin)
43.                        //如果删除的是begin元素, 重置迭代器
44.                    {
45.                        iter = common_docid.begin();
46.                    }
47.                    else
48.                    {
49.                        iter = ++ it_back;
50.                        //删除元素后重新设置迭代器
51.                    }
52.                }
53.                else
54.                {
55.                    iter++;
56.                }
57.            }
58.        }
59.    }
60.    #ifndef DEBUG //测试输出包含搜索词的docid
61.        for(set<int>::iterator iter = common_docid.begin(); iter != common_docid.end(); ++iter)
62.        {
63.            cout << "common_docid: " << *iter << endl;;
64.        }
65.    #endif
66. }
```

然后将query当做一篇doc，计算单词的权重，以及归一化权重。

根据公式，计算两篇doc（query和doc）之间的相似度：

```
1. /*
2.  * 计算相似度
3.  * 求doc与搜索关键字的相似度
4.  * 结果存入优先级队列
5.  */
6. void compute_similarity(set<int> &common_docid, map<string, double> &search_word_normalized_power, map<string, map<int, double> > &
7. {
8.     //求相似度
9.     for(set<int>::iterator iter = common_docid.begin(); iter != common_docid.end(); ++iter)
10.    {
11.        int docid = *iter;
12.        double similarity = 0;
```

```
13.         for(map<string, double>::iterator it = search_word_normalized_power.begin(); it != search_word_normalized_power.end(); it++)
14.         {
15.             similarity += it->second * word_index[it->first][docid];
16.         }
17.         Similarity sim;
18.         sim._docid = docid;
19.         sim._similarity = similarity;
20.         q.push(sim);
21.     }
22. }
```

从计算得到的docid取出整篇doc，分别取出标题和内容：

```
1.  std::string Task::get_title(const std::string &doc) //取标题
2.  {
3.      int start = doc.find("<title>") + 7;
4.      int end = doc.find("</title>");
5.      string title(doc, start, end - start);
6.      if(title[0] == '\n')
7.      {
8.          title.erase(0, 1);
9.      }
10.     return title;
11. }
```

```
1.  std::string Task::get_content(const std::string &doc) //取内容
2.  {
3.      int start = doc.find("<content>") + 9;
4.      int end = doc.find("</content>");
5.      string content(doc, start, end - start);
6.      if(content[0] == '\n')
7.      {
8.          content.erase(0, 1);
9.      }
10.     string line;
11.     istream iss(content);
12.     int count = m_out_line;
13.     string ret;
14.     while(getline(iss, line) && count > 0)
15.     {
16.         ret += line + "<br>"; //html的换行模式
17.         count--;
18.     }
19.     ret += ".....";
20.     return ret;
21. }
```

将查询结果制作成 Jason 字符串，发回客户端（前台页面）：

```
1.  /*
2.   * 将一个vector<pair<string,string> >做成json字符串
3.   *  pair中存放两个string，分别是title和content
4.   */
5.  static std::string json_string(vector<pair<string, string> > &result_pair)
6.  {
7.      Json::Value root ;
8.      Json::Value arr ;
9.      for(vector<pair<string, string> >::iterator iter = result_pair.begin(); iter != result_pair.end(); ++iter)
10.     {
11.         Json::Value elem ;
12.         elem["title"] = iter->first ;
13.         elem["summary"] = iter->second ;
14.         arr.append(elem);
15.     }
16.     root["files"]=arr ;
17.     Json::FastWriter writer ;
18.     Json::StyledWriter stlwriter ;
19.     return stlwriter.write(root);
20. }
```

第七阶段： 前台页面

index.html 使用 javascript 通过 post 方法向 php 写的 tcp_client 发送查询 query，并将接收到的查询结果（Jason字符串）解析出来后显示在页面上：

```
1.  <script>
2.  //点击search按钮，执行其中的事件，注意js代码的注释与html代码的注释的区别
3.  $("#submitButton").click(function(){
4.
5.      //取输入框的值
6.      var myWords=$("#txtSearch").val();
7.      //ajax请求，方法为post,php客户端返回的数据（echo）存在data变量中
8.      $.post("tcp_client.php",{content:myWords},function(data,status){
9.          if(status=="success")//post请求状态成功
10.         {
11.             //将收到的json字符串（data）转化为json对象，注意json字符串与json对象的区别
12.             var obj = eval("(" + data + ")");
13.             $("#result").html("");//清空result内容,用的是jquery的html()函数
14.             $.each(obj.files, function(i, item) { //遍历json对象，用的是jquery的each()方法，该json对象的格式近似于: {"files":[{"title":title_1,"summary":summary_1}]}
15.                 $("#result").append(//将遍历到的数据显示在id为result这个div里面
16.                     //根据json对象的每一个子集的键显示相应的值，哟给你的是json的语法
17.                     "<div>" + item.title + "</div>" +
18.                     "<div>" + item.summary+ "</div><hr/>");
19.             });
20.
21.     }
```



```
22.         else //post failure
23.         {
24.
25.             alert(error);
26.         }
27.
28.     }); //end post
29.
30.     });
31.
32. </script><!--  javascript 结束 -->
```

php 客户端，主要功能是，当用户提交查询时，使用 tcp 协议向服务器发送查询词，并接收服务器发回的查询结果：

```
1.  <?php
2.  $buff=$_REQUEST["content"];//采用$_REQUEST超全局数组来接收index.html页面post请求传递过来的数据
3.  //tcp client
4.  $server_Ip="127.0.0.1";//服务端ip地址，如果你的客户端与服务端不在同一台电脑，请修改该ip地址
5.  $server_Port=5080;//通信端口号
6.
7.  //设置超时时间
8.  set_time_limit(0);
9.  //创建套接字
10.  $sock= socket_create(AF_INET,SOCK_STREAM,SOL_TCP);
11.  if(!$sock)
12.  {
13.      echo "creat sock failed";
14.      exit();//创建套接字失败，结束程序
15.  }
16.
17.  socket_connect($sock,$server_Ip,$server_Port);
18.
19.  //发送数据到tcp的服务端（C语言写的）
20.  socket_send($sock,$buff,strlen($buff),0);
21.  $buff="";//清空缓冲区
22.  socket_recv($sock,$buff,1024000,0);//接收tcp_server传递过来json字符串，存在变量$buff中
23.
24.  echo trim($buff)."\n";//去掉接受到的字符串的首尾空格，返回给post请求的data
25.  //关闭套接字
26.  socket_close($sock);
27.  ?>
```

-
-
-
-
-
-
-
-

×

通知

取消 确认

内容目录

◦	▪	Cmd-Markdown Cmd Markdown 简明语法手册	1	◦	2014-05-14 项目二：miniSearch项目文档
◦	▪	Cmd-Markdown Cmd Markdown 简明语法手册	1	▪	说明
◦	▪	Linux vim 配置 YouCompleteMe	2	▪	第一阶段：搭建框架 0514
◦	▪	一句shell命令搞定代码行数统计		▪	1.目录结构
◦	▪	Linux网络编程 2014-04-04 Linux网络编程 - 基于TCP的文件传输	5	▪	2.将服务器设置为守护进程
▪	2014-04-04 Linux网络编程 - socket建立TCP连接			▪	3.封装tcp类
▪	2014-04-03 Linux网络编程 - 基于UDP的多人群聊			▪	4.搭建多进程框架
▪	2014-04-02 Linux网络编程 - socket建立UDP连接			▪	第二阶段：网页库构建 0515
▪	2014-04-01 Linux网络编程 - 大端小端、ip转换			▪	编码的转换
◦	▪	SQL 2014-06-13 SQL	1	▪	第三阶段：网页去重 0516
◦	▪	git 2014-04-18 GitHub	1	▪	第四阶段：建索引 0517
◦	▪	项目文档 2014-06-17 FTP文件服务器（TCP、HTTP）	3	▪	目标：
▪	2014-05-04 项目一：SpellCorrect项目文档			▪	第一步，根据公式求权重：
▪	2014-05-14 项目二：miniSearch项目文档			▪	第二步，根据公式，把权重归一化：
◦	▪	未分类 2014-06-11 六种排序方式 sort.c	1	▪	第三步，单词索引写入文件：
◦	▪	联系我们		▪	第五阶段：计算文本相似度 0518
◦	▪	关注开发者		▪	第六阶段：查询模块 0518
◦	▪	报告问题、建议		▪	第七阶段：前台页面

-
- 以下【标签】将用于标记这篇文稿：