

Introduction to Web Science

Assignment 5

Prof. Dr. Steffen Staab

staab@uni-koblenz.de

René Pickhardt

rpickhardt@uni-koblenz.de

Korok Sengupta

koroksengupta@uni-koblenz.de

Institute of Web Science and Technologies

Department of Computer Science

University of Koblenz-Landau

Submission until: November 30, 2016, 10:00 a.m.

Tutorial on: December 2, 2016, 12:00 p.m.

Please look at the lessons 1) **Dynamic Web Content** & 2) **How big is the Web?**

For all the assignment questions that require you to write code, make sure to include the code in the answer sheet, along with a separate python file. Where screen shots are required, please add them in the answers directly and not as separate files.

Team Name: Echo

1 Creative use of the Hypertext Transfer Protocol (10 Points)

HTTP is a request response protocol. In that spirit a client opens a TCP socket to the server, makes a request, and the server replies with a response. The server will just listen on its open socket but cannot initiate a conversation with the client on its own.

However you might have seen some interactive websites which notify you as soon as something happens on the server. An example would be Twitter. Without the need for you to refresh the page (and thus triggering a new HTTP request) they let you know that there are new tweets available for you. In this exercise we want you to make sense of that behaviour and try to reproduce it by creative use of the HTTP protocol.

Have a look at `server.py`¹ and `webclient.html` (which we provide). Extend both files in a way that after `webclient.html` is served to the user the person controlling the server has the chance to make some input at its commandline. This input should then be sent to the client and displayed automatically in the browser without requiring a reload. For that the user should not have to interact with the webpage any further.

1.1 webclient.html

```
1: <html>
2: <head>
3:     <title>Abusing the HTTP protocol - Example</title>
4: </head>
5: <body>
6:     <h1>Display data from the Server</h1>
7:     The following line changes on the servers command line
8:     input: <br>
9:     <span id="response" style="color:red">
10:         This will be replaced by messages from the server
11:     </span>
12: </body>
13: </html>
```

1.2 Hints:

- This exercise is more like a riddle. Try to focus on how TCP sockets and HTTP work and how you could make use of that to achieve the expected behaviour. Once you have an idea the programming should be straight-forward.
- The Javascript code that you need for this exercise was almost completely shown in one of the videos and is available on Wikiversity.

¹you could store the code from <http://blog.wachowicz.eu/?p=256> in a file called `server.py`

- In that sense we only ask for a "proof of concept" nothing that would be stable out in the wilde.
 - In particular, don't worry about making the server uses multithreading. It is ok to be blocking for the sake of this exercise.
- Without use of any additional libraries or AJAX framework we have been able to solve this with 19 lines of Javascript and 11 lines of Python code (we provide this information just as a way for you to estimate the complexity of the problem, don't worry about how many lines your solution uses).

Answer: Code (index.html):

```
1: <html>
2: <head>
3:     <title>Abusing the HTTP protocol - Example</title>
4: </head>
5: <body>
6:     <h1>Display data from the Server</h1>
7:     The following line changes on the servers command line
8:     input: <br>
9:     <span id="response" style="color:red">
10:         This will be replaced by messages from the server
11:     </span>
12:
13: <script>
14:
15: function httpGetAsync(theUrl, callbackFunc)
16: {
17:     var xmlHttp = new XMLHttpRequest();
18:     xmlHttp.onreadystatechange = function() {
19:         if (xmlHttp.readyState == 4 && xmlHttp.status == 200)
20:             callbackFunc(xmlHttp.responseText);
21:     }
22:
23:     xmlHttp.open("GET", theUrl, true);
24:     xmlHttp.send(null);
25: }
26:
27: function myCallBackFunc (returedDataFromServer){
28:     document.getElementById("response").innerHTML = returedDataFromServer;
29:     return true;
30: }
31:
32: var currentURL=window.location.href;
33:
34: currentURL=currentURL+"?jscall";
35:
36: setInterval(function(){
```

```
37:         httpGetAsync(currentURL, myCallBackFunc);
38:     }, 999);
39:
40:
41: </script>
42:
43:
44:
45: </body>
46: </html>
```

Code (simpleServer.py):

```
1: #!/usr/bin/python
2:
3: import socket
4: import signal
5: import time
6:
7: class Server:
8:     """ Class describing a simple HTTP server objects."""
9:
10:    def __init__(self, port = 80):
11:        """ Constructor """
12:        self.host = '' # <-- works on all available network interfaces
13:        self.port = port
14:        self.www_dir = 'www' # Directory where webpage files are stored
15:
16:    def activate_server(self):
17:        """ Attempts to acquire the socket and launch the server """
18:        self.socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
19:        self.socket.setsockopt(socket.SOL_SOCKET, socket.SO_REUSEADDR, 1)
20:
21:        try: # user provided in the __init__() port may be unavailable
22:            print("Launching HTTP server on ", self.host, ":", self.port)
23:            self.socket.bind((self.host, self.port))
24:
25:        except Exception as e:
26:            print ("Warning: Could not acquire port:", self.port, "\n")
27:            print ("I will try a higher port")
28:            # store to user provided port locally for later (in case 8080 fails)
29:            user_port = self.port
30:            self.port = 8080
31:
32:        try:
33:            print("Launching HTTP server on ", self.host, ":", self.port)
34:            self.socket.bind((self.host, self.port))
35:
```

```
36:         except Exception as e:
37:             print("ERROR: Failed to acquire sockets for ports ", user_port, " and ", server_port)
38:             print("Try running the Server in a privileged user mode.")
39:             self.shutdown()
40:             import sys
41:             sys.exit(1)
42:
43:     print ("Server successfully acquired the socket with port:", self.port)
44:
45:     self._wait_for_connections()
46:
47: def shutdown(self):
48:     """ Shut down the server """
49:     try:
50:         print("Shutting down the server")
51:         s.socket.shutdown(socket.SHUT_RDWR)
52:
53:     except Exception as e:
54:         print("Warning: could not shut down the socket. Maybe it was already closed.")
55:
56: def _gen_headers(self, code):
57:     """ Generates HTTP response Headers. Omits the first line! """
58:
59:     # determine response code
60:     h = ''
61:     if (code == 200):
62:         h = 'HTTP/1.1 200 OK\n'
63:     elif (code == 404):
64:         h = 'HTTP/1.1 404 Not Found\n'
65:
66:     # write further headers
67:     current_date = time.strftime("%a, %d %b %Y %H:%M:%S", time.localtime())
68:     h += 'Date: ' + current_date + '\n'
69:     h += 'Server: Simple-Python-HTTP-Server\n'
70:     h += 'Connection: close\n\n'
71:
72:     return h
73:
74: def _wait_for_connections(self):
75:     """ Main loop awaiting connections """
76:     while True:
77:         print ("Awaiting New connection")
78:         self.socket.listen(10)
79:
80:         conn, addr = self.socket.accept()
81:
82:
83:         print("Got connection from:", addr)
84:
```

```
85:
86:     data = conn.recv(1024)
87:     string = bytes.decode(data)
88:
89:
90:
91:     #determine request method (HEAD and GET are supported)
92:     request_method = string.split(' ')[0]
93:
94:
95:     #Checking request
96:
97:     x= string.find("jscall")
98:     if (x!=-1):
99:         userData=input("Type your message : ")
100:         conn.sendall(str.encode(userData))
101:         # socket only send the data to client after closing the connection
102:         conn.close()
103:
104:         self.activate_server()
105:         _wait_for_connections(self)
106:
107:
108:     if (request_method == 'GET') | (request_method == 'HEAD'):
109:
110:
111:         # split on space "GET /file.html" -into-> ('GET','file.html',...)
112:         file_requested = string.split(' ')
113:         file_requested = file_requested[1]
114:
115:         #Check for URL arguments. Disregard them
116:         file_requested = file_requested.split('?')[0]
117:
118:         if (file_requested == '/'):
119:             file_requested = '/index.html'
120:
121:         file_requested = self.www_dir + file_requested
122:         print ("Serving web page [" ,file_requested,"]")
123:
124:         ## Load file content
125:         try:
126:             file_handler = open(file_requested,'rb')
127:             if (request_method == 'GET'): #only read the file when GET
128:                 response_content = file_handler.read() # read file content
129:                 file_handler.close()
130:
131:                 response_headers = self._gen_headers( 200)
132:
133:         except Exception as e:
```

```
134:         print ("Warning, file not found. Serving response code 404\n", e)
135:         response_headers = self._gen_headers( 404)
136:
137:         if (request_method == 'GET'):
138:             response_content = b"<html><body><p>Error 404: File not found."
139:
140:         server_response = response_headers.encode()
141:         if (request_method == 'GET'):
142:
143:             server_response = response_content
144:
145:
146:         conn.send(server_response)
147:
148:
149:         print ("Closing connection with client")
150:         conn.close()
151:
152:     else:
153:         print("Unknown HTTP request method:", request_method)
154:
155: def graceful_shutdown(sig, dummy):
156:     """ This function shuts down the server. It's triggered
157:     by SIGINT signal """
158:     s.shutdown() #shut down the server
159:     import sys
160:     sys.exit(1)
161:
162: #####
163: # shut down on ctrl+c
164: signal.signal(signal.SIGINT, graceful_shutdown)
165:
166: print ("Starting web server")
167: s = Server(80) # construct server object
168: s.activate_server() # aquire the socket
```

2 Web Crawler (10 Points)

Your task in this exercise is to "crawl" the Simple English Wikipedia. In order to execute this task, we provide you with a mirror of the Simple English Wikipedia at 141.26.208.82.

You can start crawling from <http://141.26.208.82/articles/g/e/r/Germany.html> and you can use the `urllib` or `doGetRequest` function from the last week's assignment.

Given below is the strategy that you might adopt to complete this assignment:

1. Download <http://141.26.208.82/articles/g/e/r/Germany.html> and store the page on your file system.
2. Open the file in python and extract the local links. (Links within the same domain.)
3. Store the file to your file system.
4. Follow all the links and repeat steps 1 to 3.
5. Repeat step 4 until you have downloaded and saved all pages.

2.1 Hints:

- Before you start this exercise, please have a look at Exercise 3.
- Make really sure your crawler doesn't follow external urls to domains other than <http://141.26.208.82>. In that case you would start crawling the entire web
- Expect the crawler to run about 60 Minutes if you start it from the university network. From home your runtime will most certainly be even longer.
- It might be useful for you to have some output on the crawlers commandline depicting which URL is currently being fetched and how many URLs have been fetched so far and how many are currently on the queue.
- You can (but don't have to) make use of breadth-first search.
- It probably makes sense to take over the full paths from the pages of the Simple English Wikipedia and use the same folder structure when you save the html documents.
- You can (but you don't have to) speed up the crawler significantly if you use multithreading. However you should not use more than 10 threads in order for our mirror of Simple English Wikipedia to stay alive.

Answer: Code (crawl.py):

```
1: from urlparse import urlparse
2: from urlparse import urljoin
3: from Queue import Queue
4: from io import open
5: import socket
6: import re
7: import os
8:
9: server_ip = "http://141.26.208.82/"
10: prsurl = urlparse(server_ip)
11: global alreadydownloaded
12: global paths
13: global failed
14: global downloaded
15: downloaded = [] # local paths of the downloaded files
16: failed = []
17: paths = [] # URLs of the downloaded files
18: alreadydownloaded = []
19:
20:
21: def recvall(sock1): # receiving all data
22:     dat = ""
23:     part = None
24:     while part != "":
25:         part = sock1.recv(2046)
26:         dat += part
27:     return dat
28:
29:
30: def dohttpgetrequest(youarel): # making request and downloading file
31:     """
32:
33:     :rtype: str full path of the file
34:     """
35:     try:
36:         parsed = urlparse(youarel)
37:         path = parsed.path
38:     except:
39:         print 'URL parsing is wrong '
40:         failed.append(youarel)
41:         pass
42:     client = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
43:     port = 80
44:     crlf = "\r\n\r\n"
45:     client.connect(('141.26.208.82', port))
46:     try:
47:         request = "GET {0} HTTP/1.0{1}".format(path, crlf)
48:         client.send(request.encode())
49:     except:
```

```
50:         print 'problem in request'
51:         failed.append(path)
52:         pass
53:     data = recvall(client)
54:     body = data.split("\r\n\r\n")[1]
55:     if not os.path.exists(os.getcwd() + '/output' + os.path.dirname(path)):
56:         os.makedirs(os.getcwd() + '/output' + os.path.dirname(path))
57:     dir1 = os.getcwd() + '/output' + os.path.dirname(path)
58:     filename1 = dir1 + '/' + os.path.basename(path)
59:     if os.path.isfile(filename1):
60:         alreadydownloaded.append(filename1)
61:     else:
62:         try:
63:             filename1 = dir1 + '/' + os.path.basename(path)
64:             with open(filename1, 'wb+') as f:
65:                 f.write(body)
66:                 downloaded.append(filename1)
67:                 return filename1
68:         except IOError:
69:             print 'saving file issue ', youare1
70:             failed.append(filename1)
71:             pass
72:
73:
74: def resolve_url_path(path): # resolving the dots in the URLs
75:     """
76:
77:     :rtype: str
78:     """
79:     segments = path.split('/')
80:     segments = [segment + '/' for segment in segments[:-1]] + [segments[-1]]
81:     resolved = []
82:     for segment in segments:
83:         if segment in ('../', '..'):
84:             if resolved[1:]:
85:                 resolved.pop()
86:         elif segment not in ('./', '.'):
87:             resolved.append(segment)
88:     return ''.join(resolved)
89:
90:
91: def readfile(ff):
92:     fo = open(ff, encoding="utf8")
93:     fill = fo.read()
94:     fo.close()
95:     return fill
96:
97:
98: def findandsave(base): # finding the URL links <a tags in a Queue
```

```
99:     """
100:
101:     :rtype: Queue
102:     """
103:     internalurls = Queue()
104:     if base:
105:         try:
106:             fx = readfile(base)
107:             links = re.findall('<a href=[\'"]?([^\'" >]+)', fx, re.DOTALL)
108:             for link in links:
109:                 merge = urljoin(server_ip, link)
110:                 check = urlparse(merge).netloc
111:                 if urlparse(merge).path != "/":
112:                     if check == prsurl.netloc:
113:                         esm = urljoin(server_ip, resolve_url_path(link))
114:                         internalurls.put(esm)
115:         except:
116:             print 'something went wrong in the saveandfind function: ', base
117:     return internalurls
118:
119:
120: def downloadlist(alist): # downloading the Queues
121:     blistt = alist
122:     while not blistt.empty():
123:         x = blistt.get()
124:         if x not in paths:
125:             y = dohttpgetrequest(x)
126:             paths.append(x)
127:             z = findandsave(y)
128:             downloadlist(z)
129:
130: # downloading germany.html as a first web page
131: gerpath = "http://141.26.208.82/articles/g/e/r/Germany.html"
132: if gerpath not in paths:
133:     filename = dohttpgetrequest("http://141.26.208.82/articles/g/e/r/Germany.html")
134:     paths.append(gerpath)
135: thelist = findandsave(filename)
136: downloadlist(thelist)
137:
138: # recording data in files, NOTE they're in appending mode
139: with open('downloaded.txt', 'ab+') as fil:
140:     for down in downloaded:
141:         fil.write(down + '\n')
142: with open('failed.txt', 'ab+') as fil:
143:     for fail in failed:
144:         fil.write(fail + '\n')
145: with open('alreadydownloaded.txt', 'ab+') as fil:
146:     for already in alreadydownloaded:
147:         fil.write(already + '\n')
```

```
148: with open('paths.txt', 'ab+') as fil:  
149:     for pth in paths:  
150:         fil.write(pth + '\n')
```

3 Web Crawl Statistics (10 Points)

If you have successfully completed the first exercise of this assignment, then please provide the following details. You may have to tweak your code in the above exercise for some of the results.

3.1 Phase I

1. Total Number of *webpages* you found.
2. Total number of links that you encountered in the complete process of crawling.
3. Average and median number of links per web page.
4. Create a *histogram* showing the distribution of links on the crawled web pages. You can use a bin size of 5 and scale the axis from 0-150.

3.2 Phase II

1. For every page that you have downloaded, count the number of internal links and external links.
2. Provide a *scatter plot* with number of internal links on the X axis and number of external links on the Y axis.

Answer: Code is uploaded but buggy and not complete.

Important Notes

Submission

- Solutions have to be checked into the github repository. Use the directory name `groupname/assignment5/` in your group's repository.
- The name of the group and the names of all participating students must be listed on each submission.
- Solution format: all solutions as *one* PDF document. Programming code has to be submitted as Python code to the github repository. Upload *all* `.py` files of your program! Use UTF-8 as the file encoding. *Other encodings will not be taken into account!*
- Check that your code compiles without errors.
- Make sure your code is formatted to be easy to read.
 - Make sure you code has consistent [indentation](#).
 - Make sure you comment and document your code adequately in English.
 - Choose consistent and intuitive names for your identifiers.
- Do *not* use any accents, spaces or special characters in your filenames.

Acknowledgment

This latex template was created by Lukas Schmelzeisen for the tutorials of "Web Information Retrieval".

\LaTeX

Currently the code can only be build using [LuaLaTeX](#), so make sure you have that installed. If on Overleaf, there's an error, go to settings and change the \LaTeX engine to LuaLaTeX.