

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.

Group Name : Yankee

Candidates :

1. Sabin Bhattarai - 216203590
sbhattarai@uni-koblenz.de
2. Biplov K.C. - 216203865
biplov@uni-koblenz.de
3. Syed Salman Ali. - 216203923
salmanali@uni-koblenz.de

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 1

Please note we have provided a file index.html inside www folder during submission./

Display data from the Server

The following line changes on the servers command line input:

Hello

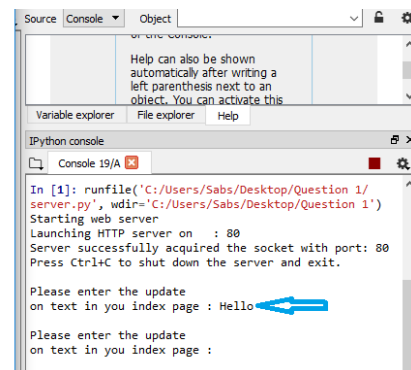


Figure 1: Shows the update when text typed in command line from server end

Below we show the Server Side code downloaded as told and edited/

```
1: #ServerSide Code
2: #!/usr/bin/python
3:
4: import socket # Networking support
5: import signal # Signal support (server shutdown on signal receive)
6: import time # Current time
7: import sys
8: from threading import Thread
9:
10:
11: inputStringFromCommandGlobal = "This will be replaced by messages from the server"
12:
13: class Server:
14:     """ Class describing a simple HTTP server objects."""
15:
16:     def __init__(self, port = 80):
17:         """ Constructor """
18:         self.host = '' # <-- works on all available network interfaces
19:         self.port = port
20:         self.www_dir = 'www' # Directory where webpage files are stored
21:
22:     def activate_server(self):
23:         """ Attempts to acquire the socket and launch the server """
24:         self.socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
25:         try: # user provided in the __init__() port may be unavailable
26:             print("Launching HTTP server on ", self.host, ":", self.port)
27:             self.socket.bind((self.host, self.port))
28:
29:         except Exception as e:
30:             print ("Warning: Could not acquire port:", self.port, "\n")
31:             print ("I will try a higher port")
32:             # store to user provided port locally for later (in case 8080 fails)
33:             user_port = self.port
34:             self.port = 8080
35:
36:         try:
37:             print("Launching HTTP server on ", self.host, ":", self.port)
38:             self.socket.bind((self.host, self.port))
39:
40:         except Exception as e:
41:             print("ERROR: Failed to acquire sockets for ports ", user_port, " \
42:                                     and 8080. ")
43:             print("Try running the Server in a privileged user mode.")
44:             self.shutdown()
45:             import sys
46:             sys.exit(1)
47:
48:     print ("Server successfully acquired the socket with port:", self.port)
```

```
49:     print ("Press Ctrl+C to shut down the server and exit.")
50:     self._wait_for_connections()
51:
52: def shutdown(self):
53:     """ Shut down the server """
54:     try:
55:         print("Shutting down the server")
56:         s.socket.shutdown(socket.SHUT_RDWR)
57:
58:     except Exception as e:
59:         print("Warning: could not shut down the socket. Maybe it was \
60:             already closed?",e)
61:
62: def _gen_headers(self, code):
63:     """ Generates HTTP response Headers. Ommits the first line! """
64:
65:     # determine response code
66:     h = ''
67:     if (code == 200):
68:         h = 'HTTP/1.1 200 OK\n'
69:     elif (code == 404):
70:         h = 'HTTP/1.1 404 Not Found\n'
71:
72:     # write further headers
73:     current_date = time.strftime("%a, %d %b %Y %H:%M:%S", time.localtime())
74:     h += 'Date: ' + current_date + '\n'
75:     h += 'Server: Simple-Python-HTTP-Server\n'
76:     # signal that the conection will be closed after complting the request
77:     h += 'Connection: close\n\n'
78:
79:     return h
80:
81: def _wait_for_connections(self):
82:     """ Main loop awaiting connections """
83:     while True:
84:         #print ("Awaiting New connection")
85:         self.socket.listen(3) # maximum number of queued connections
86:
87:         conn, addr = self.socket.accept()
88:         # conn - socket to client
89:         # addr - clients address
90:
91:         #print("Got connection from:", addr)
92:
93:         data = conn.recv(1024) #receive data from client
94:         string = bytes.decode(data) #decode it to string
95:
96:         #determine request method (HEAD and GET are supported)
97:         request_method = string.split(' ')[0]
```

```
198:         #print ("Method: ", request_method)
199:         #print ("Request body: ", string)
200:
201:         #if string[0:3] == 'GET':
202:         if (request_method == 'GET') | (request_method == 'HEAD'):
203:             #file_requested = string[4:]
204:
205:             # split on space "GET /file.html" -into-> ('GET','file.html',...)
206:             file_requested = string.split(' ')
207:             file_requested = file_requested[1] # get 2nd element
208:
209:             #Check for URL arguments. Disregard them
210:             file_requested = file_requested.split('?')[0] # disregard anything &
211:
212:             if (file_requested == '/'): # in case no file is specified by the b
213:                 file_requested = '/index.html' # load index.html by default
214:
215:             file_requested = self.www_dir + file_requested
216:             #print ("Serving web page [",file_requested,"]")
217:
218:             ## Load file content
219:             try:
220:                 file_handler = open(file_requested,'rb')
221:                 if (request_method == 'GET'): #only read the file when GET
222:                     response_content = file_handler.read() # read file content
223:                 file_handler.close()
224:
225:                 response_headers = self._gen_headers( 200)
226:
227:             except Exception as e: #in case file was not found, generate 404 pag
228:                 if(file_requested == 'www/GETINPUT' and request_method == 'GET')
229:                     response_headers = self._gen_headers( 200)
230:                     #print ("oh its a input file")
231:                     response_content = inputStringFromCommandGlobal.\
232:                                     encode('utf-8')
233:                 else:
234:                     response_headers = self._gen_headers( 404)
235:                     response_content = b"<html><body><p>Error 404: File not \
236:                                     found</p><p>Python HTTP server</p></body></html>"
237:
238:             # return headers for GET and HEAD
239:             server_response = response_headers.encode()
240:             if (request_method == 'GET'):
241:                 #return additional conten for GET only
242:                 server_response += response_content
243:
244:             conn.send(server_response)
245:             conn.close()
246:
```

```
147:
148: def graceful_shutdown(sig, dummy):
149:     """ This function shuts down the server. It's triggered
150:     by SIGINT signal """
151:     s.shutdown() #shut down the server
152:     sys.exit(1)
153:
154:
155: # shut down on ctrl+c
156: signal.signal(signal.SIGINT, graceful_shutdown)
157:
158: print ("Starting web server")
159: s = Server(80) # construct server object
160:
161: def runServer():
162:     s.activate_server() # aquire the socket
163:
164: #Starting a server thread
165: t1 = Thread(target=runServer)
166: t1.start()
167:
168: #Main Thread running for command line inputs.
169: while True:
170:     inputStringFromCommandGlobal = input("Please enter the update \
171:                                         on text in you index page : ")
```

The following is the index.html

```
1: #index.html
2: <html>
3: <head>
4:     <title>
5:         Abusing the HTTP protocol - Example
6:     </title>
7: </head>
8:
9: <body>
10:     <h1> Display data from the Server </h1>
11:     The following line changes on the servers command line input: <br>
12:     <span id="response" style="color:red">
13:         This will be replaced by messages from the server
14:     </span>
15:
16:
17:     <script language="javascript">
18:         function checkUpdates() {
19:             var XMLHttpRequestObject = false;
20:             if (window.XMLHttpRequest) {
21:                 XMLHttpRequestObject = new XMLHttpRequest();
22:             } else if (window.ActiveXObject) {
23:                 XMLHttpRequestObject = new ActiveXObject("Microsoft.XMLHTTP");
24:             }
25:             if (XMLHttpRequestObject) {
26:                 var obj = document.getElementById("response");
27:                 XMLHttpRequestObject.open("GET", "GETINPUT");
28:                 XMLHttpRequestObject.onreadystatechange = function () {
29:                     if (XMLHttpRequestObject.readyState == 4 &&
30:                         XMLHttpRequestObject.status == 200) {
31:                         obj.innerHTML = XMLHttpRequestObject.responseText;
32:                     }
33:
34:                     checkUpdates();
35:                 }
36:                 XMLHttpRequestObject.send()
37:             }
38:         }
39:         checkUpdates()
40:
41:     </script>
42:
43: </body>
44: </html>
```

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.

The following is the crawl.py

```
1: #crawl.py
2: import re
3: import socket
4: import numpy as np
5: import pandas as pd
6: import seaborn as sns
7: import matplotlib.pyplot as plt
8: from urllib.parse import urlparse
9: from urllib.parse import urljoin
10:
11: ENCODING = 'utf-8'
12: url_to_process = {'http://141.26.208.82/articles/g/e/r/Germany.html' : 0}
13: url_visited = {}
14: url_external_internal = {}
15: url_totalLink = {}
16:
17:
18: #Helper function that handles socket connection and also gives response in bytes
19: def _clientConnectAndResponse(url, timeout=10, receive_buffer=8096):
20:     parsed = urlparse(url)
21:     try:
22:         host, port = parsed.netloc.split(':')
23:     except ValueError:
24:         host, port = parsed.netloc, 80
25:
26:     try:
27:         clientSocket = socket.create_connection((host, port), timeout)
28:     except socket.timeout:
29:         raise
30:
31:
32:     method = 'GET %s HTTP/1.0\r\n\r\n' % parsed.path
33:     clientSocket.sendall(bytes(method, ENCODING))
34:
35:     response = [clientSocket.recv(receive_buffer)]
36:     while response[-1]:
37:         response.append(clientSocket.recv(receive_buffer))
38:
39:     responseInBytes = b''.join(r for r in response)
40:     return responseInBytes
41:
42:
43: # Connects to the server and receives response. writes to a file,
44: # returns filename , returns None if not 200 OK
45: def _downloadHTMLToFile(url,content_typeLength=14 ):
46:     #Connecting to Server and obtaining responses in Byte format
47:     try:
48:         responseInBytes = _clientConnectAndResponse(url)
```

```
49:     except:
50:         return None
51:
52:     headerBytes = responseInBytes[:responseInBytes.find(b'\r\n\r\n')]
53:     remainingBytes = responseInBytes[responseInBytes.find(b'\r\n\r\n'):]\
54:         .strip(b'\r\n\r\n')
55:
56:     #Here we check if HTTP responses with 200 OK
57:     if (headerBytes.find(b'\r\n')) > 0 :
58:         if headerBytes[:headerBytes.find(b'\r\n')].decode(ENCODING).lower()\
59:             != "http/1.1 200 ok":
60:             #print ("Cannot determine if HTTP 200 OK. Further Processing halts")
61:             return None
62:
63:     # Check if url has a path. This is to extract filename for writing to file.
64:     if (urlparse(url).path == ''):
65:         filename = url.split('#')[-1]
66:     else:
67:         filename = urlparse(url).path.split('/')[1]
68:
69:     # Write to a file with name given by filename variable
70:     try:
71:         with open(filename, 'wb') as htmlToFile:
72:             htmlToFile.write(remainingBytes)
73:             htmlToFile.close()
74:         return filename
75:     except IOError as e:
76:         return None
77:
78:
79: # Helper function used by parseHTMLForLinks to check if url is internal
80: def _is_internal(urlMain, urlToCheck):
81:     return urlparse(urlMain).scheme == urlparse(urlToCheck).scheme and\
82:         urlparse(urlMain).netloc == urlparse(urlToCheck).netloc
83:
84:
85: # Helper function that returns the front part of the url that is to be
86: # appended to relative urls for valid image construction
87: def _getPartToAppendToRelativeLinkURL(url):
88:     parsed = urlparse(url)
89:     try:
90:         host, port = parsed.netloc.split(':')
91:     except ValueError:
92:         host = parsed.netloc
93:     return parsed.scheme + '://' + host
94:
95:
96: # obtains a dictionary of internal and external links for a given url
97: # Done by opening the file and extracting all links. Manages unique links
```

```
198: def _parseHTMLForLinks(filename , url, content_typeLength=14):
199:     externalLinks = set()
200:     internalLinks = set()
201:     allLinks = {}
202:
203:     with open(filename , encoding = ENCODING) as html:
204:         try:
205:             content = html.read()
206:         except:
207:             allLinks['internal'] = internalLinks
208:             allLinks['external'] = externalLinks
209:             return allLinks
210:         pat = re.compile ('<a href="?\'?([^\\">]*)')
211:         linksInAHref = pat.findall(content)
212:
213:         toAppendRelativeLinkURL = _getPartToAppendToRelativeLinkURL(url)
214:         for linkInAHref in linksInAHref:
215:             if(linkInAHref.find('/://') < 0):
216:
217:                 #We have a relative link. So need to modify link
218:                 link = urljoin(toAppendRelativeLinkURL, linkInAHref)
219:
220:                 internalLinks.update([link])
221:             else:
222:                 if _is_internal(url,linkInAHref ):
223:                     internalLinks.update([linkInAHref])
224:                 else:
225:                     externalLinks.update([linkInAHref])
226:
227:             allLinks['internal'] = internalLinks
228:             allLinks['external'] = externalLinks
229:             return allLinks
230:
231: # Used by process_UrlBFS to obtain all the link for a given url. Uses
232: # two helper function to download html to file and also parse html to get links
233: def getAllLinksForThisURL(url):
234:     #making sure all links obtained inside from the url is unique
235:     links = {'internal' : {} , 'external':{}}
236:
237:     fileLocation = _downloadHTMLToFile(url)
238:     if fileLocation is not None:
239:         links = _parseHTMLForLinks(fileLocation,url)
240:
241:     return links
242:
243: # Here we obtain the URL from global variable and start the URL crawling.
244: def process_urlBFS():
245:     totalLinks = 0
246:     totalPages = 0
```

```
147: while(url_to_process):
148:     key , value = url_to_process.popitem()
149:     # Lets process the url to get more url links.
150:     # Also download the html file for this url. Returns all link info
151:     receivedInternalExternalLinks = getAllLinksForThisURL(key)
152:
153:     url_totalLink[key] = len(receivedInternalExternalLinks['internal']) + \
154:         len(receivedInternalExternalLinks['external'])
155:
156:     #This will be used to calculate pages as we take 200 ok link as pages.
157:     url_visited[key] = url_totalLink[key]
158:
159:     #key as url associated to values tuple i.e (internal, external)
160:     url_external_internal[key] = (len(receivedInternalExternalLinks['internal']
161:                                     ),len(receivedInternalExternalLinks\
162:                                     ['external']))
163:
164:     totalLinks = totalLinks + len(receivedInternalExternalLinks['internal']) \
165:         + len(receivedInternalExternalLinks\
166:               ['external'])
167:     allInternalReceivedURLLinks = receivedInternalExternalLinks['internal']
168:
169:
170:     for receivedURLLink in allInternalReceivedURLLinks:
171:         # check if contains in url_visited. Add to process queue
172:         if (receivedURLLink not in url_visited):
173:             url_to_process[receivedURLLink] = 0
174:
175: totalPages = len([x for x in url_visited.values() if x != 0])
176: #badPages404 = len([x for x in url_visited.values() if x == 0])
177: print ("\nPhase I")
178: print("\nTotal Number of webpages we found : " , totalPages)
179: print("\nTotal Number of links we encountered : " , totalLinks)
180:
181: average = totalLinks / totalPages
182: print("\nAverage number of links per web page : " , average )
183: print ("\nMedian number of links per web page : " , \
184:        calculateMedian(url_totalLink))
185: drawHistogram(url_totalLink)
186: print ("\nPhase II")
187: drawScatterDiagram(url_external_internal)
188:
189: def calculateMedian(url_totalLink):
190:     medianValue = (np.median(list(url_totalLink.values()))))
191:     return medianValue
192:
193: def drawHistogram(url_totalLink):
194:     x_1 = np.array(list(url_totalLink.values()))
195:     plt.figure(figsize=(7,5))
```

```
196:     plt.hist(x_1 , bins= range(min(x_1), 150, 5))
197:     plt.xlabel("Number of Links per page (Each bar i.e. 0-5, 5-10, ...)")
198:     plt.ylabel("Frequencies")
199:     plt.show()
200:
201: def drawScatterDiagram(url_external_internal):
202:     #preparing data for plotting
203:     tuple_internal_external = list(url_external_internal.values())
204:     internal = [x[0] for x in tuple_internal_external]
205:     external = [x[1] for x in tuple_internal_external]
206:     #print ("mero internal lsit : " , internal )
207:     #print ("mero external lsit : " , external )
208:     df1 = pd.DataFrame(index=range(len(internal)), columns=["URL" , \
209:                                                             "Number Of Internal Links",\
210:                                                             "Number Of External Links"])
211:     df1["URL"] = url_external_internal.keys()
212:     df1["Number Of Internal Links"] = internal
213:     df1["Number Of External Links"] = external
214:     print(df1)
215:     sns.lmplot(x="Number Of Internal Links", y="Number Of External Links"\
216:               , data=df1, fit_reg=False)
217:
218: def mainFunction():
219:     process_urlBFS()
220:     print('done')
221:
222:
223:
224: if __name__ == "__main__":
225:     mainFunction()
```

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.

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.