# Introduction to Web Science

**Assignment 5**

Prof. Dr. Steffen Staab        René Pickhardt

staab@uni-koblenz.de        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: India
Group Member : Jalpa Patel , Jasvinder Kaur , Amani Gaddamedi

# 1 Creative use of the Hyptertext 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 servered to the user the person controlling the server has the chance to make some input at its commandline. This input should then be send 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 focuse 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.

---

[1]you 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. server.py

```
 1: import socket  # Networking support
 2: import signal  # Signal support (server shutdown on signal receive)
 3: import time    # Current time
 4:
 5: message_before = 'Intro Web Science';
 6: message_after = ""
 7:
 8: class Server:
 9:  """ Class describing a simple HTTP server objects."""
10:
11:  def __init__(self, port = 8080):
12:      """ Constructor """
13:      self.host = ''   # <-- works on all avaivable network interfaces
14:      self.port = port
15:      self.www_dir = '/home/jp/ass5' # Directory where webpage files are store
16:
17:
18:  def activate_server(self):
19:      """ Attempts to aquire the socket and launch the server """
20:      self.socket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
21:      try: # user provided in the __init__() port may be unavaivable
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 aquite port:",self.port,"\n")
27:          print ("I will try a higher port")
28:          # store to user provideed 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,
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:        print ("Press Ctrl+C to shut down the server and exit.")
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
55:
56:    def _gen_headers(self,  code):
57:        """ Generates HTTP response Headers. Ommits 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'  # signal that the conection wil be closed
71:
72:        return h
73:
74:    def long_pooling(self, conn):
75:        global message_after, message_before
76:        message_after = input("Your Message Here ")
77:        if not message_after == '' and not message_after == message_before:
78:            message_before = message_after
79:        else:
80:            message_after = message_before
81:
82:        response_headers = self._gen_headers( 200)
83:        response_content = response_headers.encode()
```

```
 84:        response_content += message_after.encode()
 85:        print(response_content)
 86:        conn.send(response_content)
 87:
 88:  def _wait_for_connections(self):
 89:        """ Main loop awaiting connections """
 90:        while True:
 91:            print ("Awaiting New connection")
 92:            self.socket.listen(3) # maximum number of queued connections
 93:
 94:            conn, addr = self.socket.accept()
 95:            # conn - socket to client
 96:            # addr - clients address
 97:
 98:            print("Got connection from:", addr)
 99:
100:            data = conn.recv(1024) #receive data from client
101:            string = bytes.decode(data) #decode it to string
102:
103:            #determine request method  (HEAD and GET are supported)
104:            request_method = string.split(' ')[0]
105:            print ("Method: ", request_method)
106:            print ("Request body: ", string)
107:
108:            #if string[0:3] == 'GET':
109:            if (request_method == 'GET') | (request_method == 'HEAD'):
110:                #file_requested = string[4:]
111:
112:                # split on space "GET /file.html" -into-> ('GET','file.html',...
113:                file_requested = string.split(' ')
114:                file_requested = file_requested[1] # get 2nd element
115:
116:                #Check for URL arguments. Disregard them
117:                file_requested = file_requested.split('?')[0]  # disregard anyth
118:
119:                if file_requested == '/server_message':
120:                    self.long_pooling(conn)
121:                else:
122:                    if (file_requested == '/'):  # in case no file is specified
123:                        file_requested = '/webclient.html' # load index.html by
124:
125:                    file_requested = self.www_dir + file_requested
126:                    print ("Serving web page [",file_requested,"]")
127:
128:                    ## Load file content
129:                    try:
130:                        file_handler = open(file_requested,'rb')
131:                        if (request_method == 'GET'):  #only read the file when
132:                            response_content = file_handler.read() # read file c
```

```
133:                        file_handler.close()
134:
135:                        response_headers = self._gen_headers( 200)
136:
137:                    except Exception as e: #in case file was not found, generate
138:                        print ("Warning, file not found. Serving response code 4
139:                        response_headers = self._gen_headers( 404)
140:
141:                        if (request_method == 'GET'):
142:                            response_content = b"<html><body><p>Error 404: File n
143:
144:                    server_response =  response_headers.encode() # return header
145:                    if (request_method == 'GET'):
146:                        server_response +=  response_content  # return additiona
147:                    conn.send(server_response)
148:                print("Closing connection with client")
149:                conn.close()
150:
151:            else:
152:                print("Unknown HTTP request method:", request_method)
153:
154: def graceful_shutdown(sig, dummy):
155:     """ This function shuts down the server. It's triggered
156:     by SIGINT signal """
157:     s.shutdown() #shut down the server
158:     import sys
159:     sys.exit(1)
160:
161: ###########################################################
162: # shut down on ctrl+c
163: signal.signal(signal.SIGINT, graceful_shutdown)
164:
165: print ("Starting web server")
166: s = Server(8080)  # construct server object
167: s.activate_server() # aquire the socket
```

2. 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 without page reload</h1>
 7:         New message from the server : <br>
 8:         <span id="response_message" style="color:blue">
 9:                This will be replaced by messages from the server
10:         </span>
11:         <script>
12:           var response = document.getElementById('response_message')
```

```
13:            function reqListener () {
14:               response.innerHTML = this.responseText;
15:            }
16:            function createRequest() {
17:              var oReq = new XMLHttpRequest();
18:              oReq.addEventListener("load", reqListener);
19:              oReq.open("GET", "http://localhost:8080/server_message");
20:              oReq.send();
21:            }
22:            setInterval(createRequest, 1000);
23:          </script>
24:
25: </body>
26: </html>
```

3. output:



**Figure 1:** Server create socket



**Figure 2:** Client Connect

**Figure 3:** Server accept request



**Figure 4:** Command line argument



**Figure 5:** User get Notification without page reload

# 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.

1. We are using the code from the last assignment to crawle the webpage.

2. we are considering the webpages that returns 200 OK status.

**Answer:**

1. crawler.py

```
 1: import http_client
 2: import csv
 3: from collections import deque
 4: import sys
 5: import re
 6: #import downloadEverything
 7: import http_client
 8:
 9: queue = deque()
10: traversed_set = set()
11: external_urls = {}
12: internals_urls = {}
13: domain =  "http://141.26.208.82/"
14:
15: def find_all_links(file_name):
16:   with open(file_name, 'r') as f:
17:     html = f.read()
18:   href_pattern = 'href=[\'"]?([^\'" >]+)'
19:   return re.findall(href_pattern, html)
20:
21: def if_relative_url(url):
22:   parse_url = http_client.parse_url(url)
23:   if(parse_url.netloc):
24:     return False
25:   else:
26:     return True
27:
28: def make_absolute_url(relative_url):
29:   relative_url = http_client.parse_url(relative_url).path
30:   return "".join([domain,relative_url] )
31:
32: def add_value_to_queue(value):
33:   if value not in traversed_set:
34:     queue.append(value)
35:     traversed_set.add(value)
36:
37: def process_queue():
38:   return queue.popleft()
39:
40: def crawl(start_url):
41:   add_value_to_queue(start_url)
42:   i = 0
43:   while(len(queue)):
44:     if(i%25 == 0 ):
45:         print('crawling\t' + str(i) + ' th page')
```

```
46:     i = i +1
47:
48:     url_to_be_crawled = process_queue()
49:     print("\ncarwaling url\n" + url_to_be_crawled)
50:     output = http_client.call_http_server(url_to_be_crawled, False)
51:     if output[0] == True:
52:       urls = find_all_links(output[1])
53:       for url in urls:
54:         if(if_relative_url(url)):
55:           absolute_url = make_absolute_url(url)
56:           url_list = internals_urls.get(url_to_be_crawled, [])
57:           add_value_to_queue(absolute_url)
58:           url_list.append(absolute_url)
59:           internals_urls[url_to_be_crawled] = url_list
60:         else:
61:           url_list = external_urls.get(url_to_be_crawled, [])
62:           url_list.append(url)
63:           external_urls[url_to_be_crawled] = url_list
64:   write = csv.writer(open("internal_urls.csv", "w"))
65:   for key, val in internals_urls.items():
66:     write.writerow([key, val])
67:   write = csv.writer(open("external_urls.csv", "w"))
68:   for key, val in internals_urls.items():
69:     write.writerow([key, val])
70:
71:
72: if __name__ == '__main__':
73:   start_url = "http://141.26.208.82/articles/g/e/r/Germany.html"
74:   crawl(start_url)
```

2. httpclient.py

```
 1: import socket
 2: from urllib.parse import urlparse
 3: import sys
 4: import os
 5: __author__ = "jasvinder"
 6:
 7:
 8: # creates an INET, STREAMING socket
 9: # input:
10: #   url:String - takes the valid url to make tcp connection
11: #   port: Integer - take the port number to connect
12: # output:
13: #   TCPSocket - returns a tcp socket object for further communication.
14: def tcp_connection(url, port):
15:   clientSocket = socket.socket(socket.AF_INET, socket.SOCK_STREAM)
16:   clientSocket.connect((url, port))
17:   return clientSocket
18:
```

```
19: # Parses the URL into six components.i.e scheme, network location, path, para
20: # input:
21: #   url_name:String - takes the valid url passed
22: # output:
23: #   String - parsed Url into 6 components
24: def parse_url(url_name):
25:   return urlparse(url_name)
26:
27: # User request for sending user input URL as HTTP request ending with two emp
28: # input:
29: #   url_name:String - takes the valid url passed
30: #   method_type:String - takes Http method type,i.e generally, POST or GET
31: # output:
32: #   get_request:String - returns the valid Http Request according to method t
33: def encode_http_request(path, method_type):
34:   protocol_type = 'HTTP/1.0'
35:   get_request = ''.join([method_type,' ', path,' ',protocol_type, '\r','\n','
36:   return get_request
37:
38: # Sends the passed http_request request over passed tcp socket object.
39: # input:
40: #   tcp_connection:TCP Socket object - takes the object of the TCP connection
41: #   http_request:String - takes the valid Http Request according to method ty
42: # output:
43: #   Bytes: returns the Http request made in bytes to send it further for comm
44: def send_http_request(tcp_connection, http_request):
45:   return tcp_connection.send(str.encode(http_request))
46:
47: # Reads the Http header of the response from tcp_connection socket.
48: # input:
49: #   tcp_connection:TCP Socket object - takes the object of the TCP connection
50: #output:
51: #   http_headers: String - gives the HTTP header of the response
52: def read_http_header(tcp_connection):
53:   http_headers = ['HTTP/1.1 200 OK']# initialising headers too 22 as it has a
54:   while(True):
55:     response = tcp_connection.recv(1).decode('utf-8')
56:     if(response == '\r'):
57:       http_headers.append(response)
58:       response = tcp_connection.recv(3).decode('utf-8')
59:       if(response == '\n\r\n'):
60:         http_headers.append(response)
61:         break
62:       else:
63:         http_headers.append(response)
64:
65:     else:
66:       http_headers.append(response)
67:   return ''.join(http_headers)
```

```
68:
69: # Checks whether the response is 200 or not from the tcp socket.
70: # input:
71: #    tcp_connection:TCP Socket object - takes the object of the TCP connection
72: #output:
73: #     Boolean : True for 200 False for all other http status code.
74: def http_status_200(tcp_connection):
75:   http_status = tcp_connection.recv(15).decode('utf-8')
76:   if(http_status.startswith('HTTP/1.1 200 OK')):
77:     return True
78:   return False
79:
80: # Reads the body of the Http response as binary
81: # input:
82: #    tcp_connection:TCP Socket object - takes the object of the TCP connection
83: #output:
84: #    list - gives the body of the Http response, as binary data
85: def read_http_body_as_binary(tcp_connection):
86:   http_body = []
87:   response = tcp_connection.recv(1024)
88:   while(response):
89:     http_body.append(response)
90:     response = tcp_connection.recv(1024)
91:   return http_body
92:
93: # Decodes http header
94: # input:
95: #    http_header: Http header string.
96: #output:
97: #     list - list of http header strings.
98: def decode_http_header(http_header):
99:   return http_header.strip().split('\r\n')
100:
101: # Reads the decoded http header
102: # input:
103: #    decoded_http_header: list of decoded http header
104: #output:
105: #     'String: tells whether type is text in binary
106: def find_content_type(decoded_http_header):
107:   #print(decoded_http_header)
108:   for value in decoded_http_header:
109:     if(value.startswith('Content-Type:')):
110:       if(value.startswith('Content-Type: text/html')):
111:         return 'text'
112:       else:
113:         return 'binary'
114:
115: # Reads the string and the file name and saves the string on the file.
116: # input:
```

```
117: #    string_data_to_be_written:String - String of text to be written on the fi
118: #    file_name: String- file name.
119: #output:
120: #     None
121: def save_text_data(string_data_to_be_written, file_name):
122:   with open(file_name, 'w+') as file:
123:     file.write(string_data_to_be_written)
124:
125: # Reads the string and the file name and saves the string on the file.
126: # input:
127: #    binary_data_array_to_be_written:List - List of Byte to be written on the
128: #    file_name: String- file name.
129: #output:
130: #     None
131: def save_binary_data(binary_data_array_to_be_written, file_name):
132:     #binary_data_to_be_written = str.encode(string_data_to_be_written)
133:     with open(file_name, 'bw') as file:
134:        for value in binary_data_array_to_be_written:
135:          file.write(value)
136:
137:
138:
139: # Reads the string and the file name and saves the string on the file.
140: # input:
141: #    url_name:String - Url of the resource to be downloaded.
142: #    saveheader: Boolean- default is True. This is a flag for stating whether
143: #output:
144: #     None
145: def call_http_server(url_name, saveheader = True):
146:
147:   parsed_url = parse_url(url_name)
148:   tcp_socket = tcp_connection(parsed_url.netloc, 80)
149:   http_request = encode_http_request(parsed_url.path, 'GET')
150:   send_http_request(tcp_socket, http_request)
151:   if(http_status_200(tcp_socket)):
152:     http_header = read_http_header(tcp_socket)
153:     http_data = read_http_body_as_binary(tcp_socket)
154:     decoded_message = decode_http_header(http_header)
155:     file_name = parsed_url.path.split('/')[-1]
156:     if saveheader:
157:       header_file_name = file_name +  '.header'
158:       save_text_data(http_header, header_file_name)
159:     if(find_content_type(decoded_message) == 'text'):
160:       text_file_name =  file_name + '.html'
161:       save_binary_data(http_data, text_file_name)
162:       tcp_socket.close()
163:       return (True, text_file_name)
164:     elif(find_content_type(decoded_message) == 'binary'):
165:       save_binary_data(http_data, file_name)
```

```
166:        tcp_socket.close()
167:        return (False, file_name)
168:    else:
169:        tcp_socket.close()
170:        return (False, 404)
171: if __name__ == "__main__":
172:    call_http_server(sys.argv[1])
```

3. output:

- 



**Figure 6:** crawl output

# 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.

1. for pages encountered we are considering all pages no duplicate removed but for url we are only considering unique 200 pages.

2. All statistical data is shared sing two csv files.They are pushed so while running they should be in home directory.

**Answer:**

1. statistics.py

```
1: # -*- coding: utf-8 -*-
2: import matplotlib.pyplot as plt
3: import numpy as np
4: import matplotlib.pyplot as hist
5: import csv
6: """
7: Created on Tue Nov 29 19:17:08 2016
8:
9: @author: Amani
10: """
11:
12: internal_urls = {}
```

```
13: external_urls = {}
14: def read_csv_as_dictionary(file_name):
15:     read_dictionary = {}
16:     for key, val in csv.reader(open(file_name)):
17:         read_dictionary[key] = val
18:     return read_dictionary
19:
20: def counting_dictionary(url_dictionary):
21:     count_urls = dict()
22:     for key, value in url_dictionary.items():
23:         count_urls[key] = len(value)
24:     return count_urls
25: # sum of values and number of keys
26: def sum_link(count_dict):
27:     sum_of_values = 0
28:     number_of_keys = 0
29:     for key, value in count_dict.items():
30:         number_of_keys += 1
31:         sum_of_values += value
32:     return (number_of_keys, sum_of_values)
33:
34:
35: # page_link gives the urls with number of links
36: def mean(page_links):
37:    values = sum_link(page_links)
38:    mean_links = values[1] / values[0]
39:    return mean_links
40:
41:
42: #function for caluculating median
43: def median(page_link):
44:         sort_list= sorted(page_link.values())
45:         length = len(sort_list)
46:         median_pos = length // 2
47:
48:         if length % 2 == 0:
49:             median = (sort_list[median_pos - 1] +
50:                        sort_list[ ( median_pos) ]) / 2
51:             return median
52:         else:
53:             median = sort_list[median_pos]
54:             return median
55: #Number of internal links
56: def add_data(data_dict1, data_dict2):
57:     added_dict  = {}
58:     for key, values in data_dict1.items():
59:         added_dict[key] = data_dict1[key] + data_dict2[key]
60:     return added_dict
61: def scaterring_graph(array_x,array_y):
```

```
62:        plt.plot(array_x,array_y, 'bo')
63:        plt.title('Scatter graph to number of internal links and external links v
64:        plt.xlabel('internal links')
65:        plt.ylabel('external links')
66:        return plt
67:
68: if __name__ == '__main__':
69:
70: #Reading the data from csv files assigned to string variables
71:        internal_link_file = 'internal_urls.csv'
72:        external_link_file = 'external_urls.csv'
73: #Assign the variables to read csv fiels
74:        internal_urls = read_csv_as_dictionary(internal_link_file)
75:        external_urls = read_csv_as_dictionary(external_link_file)
76:        external_urls_count = counting_dictionary(external_urls)
77:        internal_links_count = counting_dictionary(internal_urls)
78:        sum_of_urls = add_data(internal_links_count, external_urls_count)
79: #Printing the number of internal links
80:        print('total number of url encountered')
81:        total_urls_encountered = len(sum_of_urls.keys())
82:        print(sum_link(sum_of_urls)[0])
83: # Printing the number of externl links
84:
85: #Printing total number of urls
86:
87:
88:
89:
90:
91:        print('toal number of links encountered')
92:        print(sum_link(sum_of_urls)[1])
93:
94: #Printing the number mean and median of the links per webpage
95:        internal_link_x = [value for (key, value) in sorted(internal_links_coun
96:        external_link_y = [value for (key, value) in sorted(external_urls_count
97:        plt1 = scaterring_graph(internal_link_x, external_link_y)
98:        plt1.show()
99:        print('average number of url')
100:       print(mean(sum_of_urls))
101:       print('median number of url')
102:       print(median(sum_of_urls))
103: #Ploting the scatterring graph
104:
105:
106: #printing the histogram
107:       data = list(sum_of_urls.values())
108:       d = plt.hist(data, bins = 30)
109:       plt.show()
```
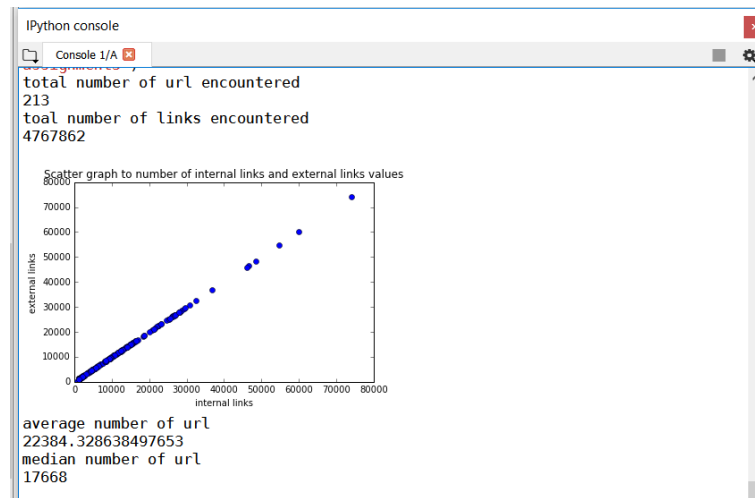
2. output:



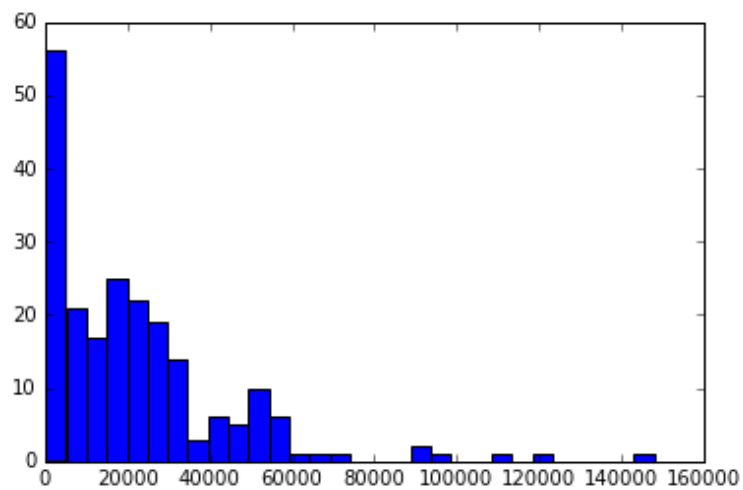**Figure 7:** Output for Statistical output



**Figure 8:** Output for statistical output

## 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 LaTeXengine to `LuaLaTeX`.