

```
1 # Introduction to Web Scraping with Python
2 # NaLette Brodnax nbrodnax@indiana.edu
3 # February 5, 2016
4
5 # ACCESS #
6 # Import all the libraries that you need
7 import requests
8 import bs4
9 import csv
10
11 webpage = 'http://www.amstat.org/publications/jse/
12 jse_data_archive.htm'
13
14
15 # PARSE #
16 soup = bs4.BeautifulSoup(server_response.text)
17
18 # create a list of dictionaries (one dict for each link)
19 link_info_list = []
20 for tag in soup.find_all('a'):
21     link = tag['href']
22     name = tag.text
23     # print(name)
24     if name[-3:] == 'txt':
25         link_info_list.append({'link': link, 'name': name})
26
27
28 # TRANSFORM #
29 # add a new category to each dictionary to categorize the
30 # file as either
31 # data or documentation
32 host = 'http://www.amstat.org/publications/jse/'
33 for dataset in link_info_list[:3]:
34     url = host + dataset['link']
35     data_response = requests.get(url)
36     if data_response.text[:5] == 'NAME:':
37         dataset['type'] = 'doc'
38     else:
39         dataset['type'] = 'dat'
40
41 # STORE #
42 def download_to_txt(file_name, data):
43     with open(file_name, 'w') as txtfile:
44         txtfile.writelines(data)
45
46
47 def strip_extension(file_name):
48     i = -1 # start at the end of the filename
```

```
49     while i > -len(file_name):
50         if file_name[i] == '.':
51             break # stop when you get to a period
52         else:
53             i -= 1 # this is the same as i = i - 1
54     return file_name[:i]
55
56 # store individual data files as text files
57 for dataset in link_info_list[:3]:
58     url = host + dataset['link']
59     data_response = requests.get(url)
60     description = strip_extension(dataset['name'])
61     filename = description + '_' + dataset['type'] + '.txt'
62     download_to_txt(filename, data_response.text)
63
64 # store list of links as csv file
65 with open('data_links.csv', 'w') as csvfile:
66     fieldnames = ['link', 'name', 'type']
67     writer = csv.DictWriter(csvfile, fieldnames)
68     writer.writeheader()
69     for link in link_info_list:
70         writer.writerow(link)
71     print('Links added to file: ' + str(len(link_info_list))
72 ))
```