

# Statistical Computing

## Lesson 4

Gabriele Rovigatti

Bank of Italy

March 19, 2021



# This Course

- Python basics
  - General Framework
  - Syntax
  - Modules / Libraries
  - Functions
- Webscraping
  - Data Storage
  - HTML Syntax
  - Data Collection
  - Requests Package
    - GET requests
    - Save Scraped Data
    - WikiScape
  - Webdriver Package
- Machine Learning?

# Today

- Python basics
  - General Framework
  - Syntax
  - Modules / Libraries
  - Functions
- **Webscraping**
  - Data Storage
  - HTML Syntax
  - Data Collection
  - **Requests Package**
    - **GET requests**
    - **Save Scraped Data**
    - **WikiScape**
  - Webdriver Package
- Machine Learning?

# HTML basics

HyperText Markup Language is the language most web pages are built with. It is **not** a programming language - it is a markup language that tells a browser how to layout content.

It consists of elements called tags. A simple example reads

```
<html>
```

```
<head>
```

```
</head>
```

```
<body>
```

```
<p>
```

This generates a line of text

```
<a href="https://www.python.org">This one generates a link to www.python.org</a>
```

```
</p>
```

```
</body>
```

```
</html>
```

# HTML-structured pages: what you see

The screenshot shows a web browser window displaying the page "Cosponsors - H.R.6729 - 112th Congress (2011-2012)". The URL is <https://www.congress.gov/bill/112th-congress/house-bill/6729/cosponsors?r=1>. The page features a navigation bar with tabs for Summary (1), Text (1), Actions (4), Titles (1), Amendments (0), **Cosponsors (9)**, Committees (2), and Related Bills (2). A sidebar on the left includes links for "BACK TO RESULTS" and "Party" (Check all, ☐ Republican). The main content area displays the title "Cosponsors: H.R.6729 — 112th Congress (2011-2012)" and a table of cosponsors. The table has columns for "Party", "Cosponsor", and "Date Cosponsored". The cosponsors listed are: Rep. Hultgren, Randy (R-IL-14)\*, Rep. Bartlett, Roscoe G. (R-MD-6)\*, Rep. McKinley, David B. (R-WV-1)\*, Rep. Bilbray, Brian P. (R-CA-50)\*, Rep. Lance, Leonard (R-NJ-7)\*, Rep. Bartlett, Lou (R-PA-11)\*, Rep. Olson, Pete (R-TX-22)\*, Rep. Thompson, Glenn (R-PA-6)\*, and Rep. Chabot, Steve (R-OH-1)\*. All dates are 01/01/2013. A "Show less" link is at the bottom of the table. The Windows taskbar at the bottom shows the time as 17:07 on 25/02/2017.

Party	Cosponsor	Date Cosponsored
<input type="checkbox"/> Republican	Rep. Hultgren, Randy (R-IL-14)*	01/01/2013
	Rep. Bartlett, Roscoe G. (R-MD-6)*	01/01/2013
	Rep. McKinley, David B. (R-WV-1)*	01/01/2013
	Rep. Bilbray, Brian P. (R-CA-50)*	01/01/2013
	Rep. Lance, Leonard (R-NJ-7)*	01/01/2013
	Rep. Bartlett, Lou (R-PA-11)*	01/01/2013
	Rep. Olson, Pete (R-TX-22)*	01/01/2013
	Rep. Thompson, Glenn (R-PA-6)*	01/01/2013
	Rep. Chabot, Steve (R-OH-1)*	01/01/2013

# HTML-structured pages: what's beneath



```
446         <tr>
447             <th scope="col" class="actions">Cosponsor</th><th scope="col" class="date">Date Cosponsored</th></tr>
448         </thead>
449         <tbody><tr>
450             <td class="actions">
451                 <a target="_blank" href="https://www.congress.gov/member/andy-hultgren/H001059">Rep. Hultgren, Randy [R-IL-14] 6#42</a>
452             </td>
453             <td class="date">01/01/2013</td>
454         </tr>
455         <tr>
456             <td class="actions">
457                 <a target="_blank" href="https://www.congress.gov/member/roscoe-bartlett/B000208">Rep. Bartlett, Roscoe G. [R-MD-6] 6#42</a>
458             </td>
459             <td class="date">01/01/2013</td>
460         </tr>
461         <tr>
462             <td class="actions">
463                 <a target="_blank" href="https://www.congress.gov/member/david-mckinley/M001180">Rep. McKinley, David B. [R-WV-1] 6#42</a>
464             </td>
465             <td class="date">01/01/2013</td>
466         </tr>
467         <tr>
468             <td class="actions">
469                 <a target="_blank" href="https://www.congress.gov/member/brian-bilbray/B000461">Rep. Bilbray, Brian P. [R-CA-50] 6#42</a>
470             </td>
471             <td class="date">01/01/2013</td>
472         </tr>
473         <tr>
474             <td class="actions">
475                 <a target="_blank" href="https://www.congress.gov/member/leonard-lance/L000567">Rep. Lance, Leonard [R-NJ-7] 6#42</a>
476             </td>
477             <td class="date">01/01/2013</td>
478         </tr>
479         <tr>
480             <td class="actions">
481                 <a target="_blank" href="https://www.congress.gov/member/lou-barletta/B001269">Rep. Barletta, Lou [R-PA-11] 6#42</a>
482             </td>
483             <td class="date">01/01/2013</td>
484         </tr>
```

# Getting to know your enemy - “Inspect” Webpages

The screenshot shows a web browser window displaying legislative search results for H.R. 6729. The browser's address bar shows the URL: `https://www.congress.gov/search?q=[\"source\": \"legislation\", \"congress\": 112, \"type\": \"bills\", \"ch...`. The search results page includes a refined filter for 'Legislation' and '112 (2011-2012)'. The search results show 'H.R. 6729' as a bill introduced in the 112th Congress (2011-2012). The description states: 'To save at least \$10,000,000,000 by consolidating some duplicative and overlapping Government programs.' The sponsor is listed as 'Rep. Guberson, John Abney (R-TX-7)' and the committee is 'House - Oversight and Government Reform, Appropriations'. The latest action is '01/01/2013 Referred to House Appropriations (All Actions)'. The tracker is 'Introduced'.

The browser's developer tools are open, showing the 'Analisi pagina' (Page Analysis) tab. The DOM tree on the left shows the structure of the search results, with the following HTML structure highlighted:

```
<div class=\"saved-search-wrapper\"></div>
<div id=\"searchTune\" class=\"basic-search-tune\"></div>
<div class=\"basic-search-results-inner-wrapper\">
  <div id=\"main\" class=\"col2_lg basic-search-results nav-on role=\"main\">
    <ol class=\"basic-search-results-lists expanded-view\" start=\"1\">
      <li class=\"expanded\" style=\"display: block;\">
        <div></div>
        1.
        <span class=\"result-heading\">
          <a href=\"https://www.congress.gov/bills/112th-congress/house-bill/6729rc=1\">H.R. 6729</a>
          - 112th Congress (2011-2012)
          :after
        </span>
        <span class=\"result-title\"></span>
        <span class=\"result-item\"></span>
        <span class=\"result-item\"></span>
        <span class=\"result-item\"></span>
        <span class=\"result-item\"></span>
        <span class=\"result-item result-tracker\"></span>
      </li>
    </ol>
  </div>
</div>
```

The right pane of the developer tools shows the 'Regole' (Rules) tab, displaying the 'Anteprima testo' (Text Preview) for the selected element. The preview shows the text 'Abc' in 'Arial' font, with a note 'Utilizzato come: \"Arial\"'.

# Last year's example: Wikipedia

The screenshot shows the Wikipedia page for 'Stazione di Tavazzano'. The page is in Italian and includes a sidebar with navigation links, a main content area with an introduction and a list of references, and a right-hand sidebar with a photo of the station and a table of location and characteristics. The page is viewed in a web browser with the address bar showing the URL 'https://it.wikipedia.org/wiki/Stazione\_di\_Tavazzano'.

Stazione di Tavazzano - Wikipedia - Modifica wikitesto

Voce [Discussione](#) [Leggi](#) [Modifica](#) [Modifica wikitesto](#) [Cronologia](#)

## Stazione di Tavazzano

Da Wikipedia, l'enciclopedia libera.

Coordinate: 45°19′35.4″N 9°24′10.51″E﻿ / ﻿(Mappa)

La **stazione di Tavazzano** è una stazione ferroviaria posta sulla linea **Milano-Bologna**.

Serve il comune di **Tavazzano con Villavesco**, ed è molto utilizzata anche dai passeggeri provenienti dai comuni confinanti **Lodi Vecchio**<sup>[1]</sup> e **Montanaso Lombardo**.

**Indice** [nascondi]

- 1 **Storia**
- 2 **Strutture ed impianti**
- 3 **Movimento**
- 4 **Progetti**
- 5 **Note**
- 6 **Altri progetti**

### Storia

[ modifica | modifica wikitesto ]

La stazione fu aperta nel **1861** con l'attivazione della **ferrovia Milano-Bologna**<sup>[2]</sup>, divenuta successivamente parte del grande itinerario dorsale italiano, da **Milano** a **Roma**. Fu costruita in aperta campagna<sup>[3]</sup>, al solo scopo di spezzare con una stazione intermedia la lunga tratta da **Melegnano** a **Lodi**.

Intorno alla stazione si sviluppò il centro abitato di **Tavazzano**, divenuto in seguito preminente sul territorio circostante, anche a scapito del **capoluogo comunale** (**Villavesco**).

Il 29 maggio **2005** divenne punto di diramazione dell'*Interconnessione di Tavazzano*, che collega la

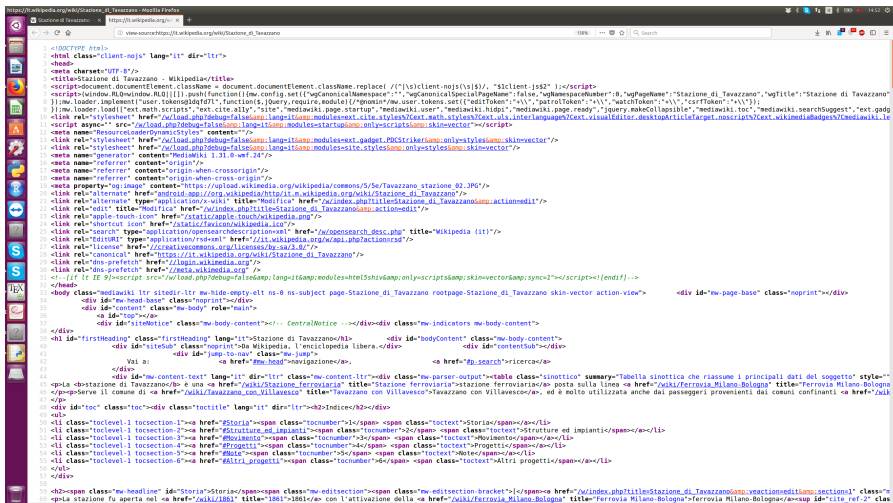
### Tavazzano stazione ferroviaria

Il fabbricato viaggiatori

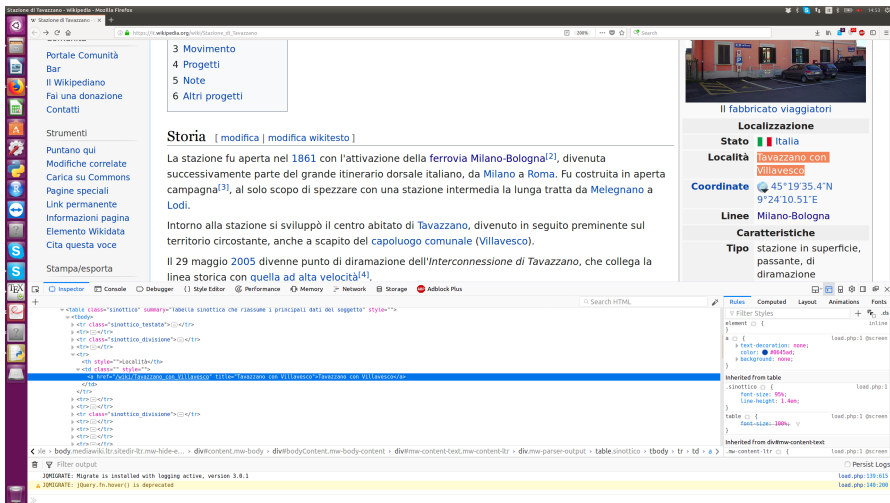
Localizzazione
<b>Stato</b> <span><span></span><span> </span></span> Italia
<b>Località</b> <span><span></span><span> </span></span> Tavazzano con Villavesco
<b>Coordinate</b> <span><span><span><span><span>45°19′35.4″N</span> <span>9°24′10.51″E</span></span></span><span><span>﻿</span> / <span>﻿</span></span><span><span></span></span></span></span>
<b>Linee</b> <span><span></span><span> </span></span> <b>Milano-Bologna</b>
Caratteristiche
<b>Tipo</b> <span><span></span><span> </span></span> stazione in superficie, passante, di



## Last year's example: Wikipedia's HTML Structure



## Last year's example: Inspect Wikipedia



# Last week's assignment - **Mandatory**

- Write a function called **CipCiop** that takes a number and
  - If the number is divisible by 3, it should return “Cip”;
  - If it is divisible by 5, it should return “Ciop”;
  - If it is divisible by both 3 and 5, it should return “CipCiop”;
  - Otherwise, it should return the original number.
- Write a function that takes as input an integer (called *limit*) and prints all the prime numbers between 0 and limit.
- Program a Rock-Paper-Scissors game, with the input() method (twice!), comparing the plays and printing the winner. (Yes, it is the same exercise for which you have the solutions, I will check your code. Try not to copy!)

# Last week's assignment - **Complimentary**

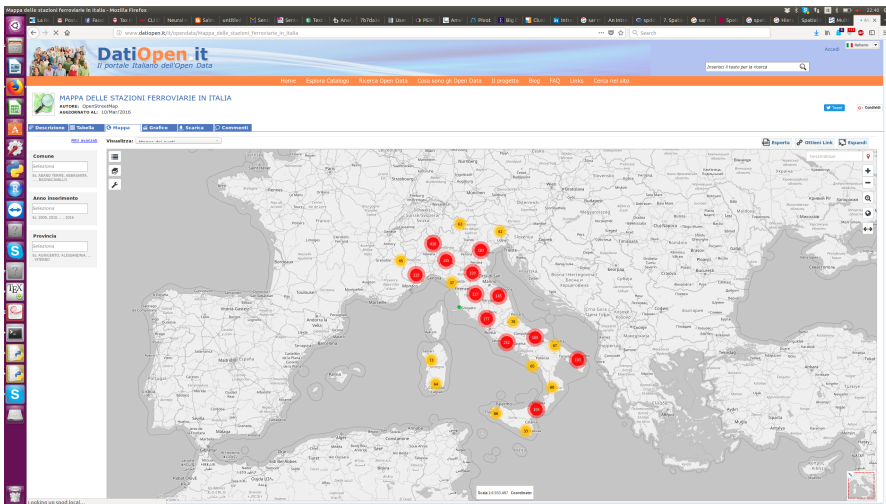
- Write a program that takes two lists as input and returns the common elements between the two lists. Make sure your program works on two lists of different sizes. [Solution](#)
- Write a one-line python code that takes a numeric list as input and returns the even elements of the list; in **one** line. [Solution](#)

# Set up a project

- Final Goal: scrape info on **all** Italian railstations
  - Station Name
  - Station Municipality
  - Opening Date
  - Closing Date - if applicable
  - Municipal Istat Code
  - Railway Line(s)
- Requirements:
  - List of all Italian train stations
  - Relative links to Wikipedia pages

# Train stations data

► [http://www.datiopen.it/it/opendata/Mappa\\_delle\\_stazioni\\_ferrovie\\_in\\_Italia](http://www.datiopen.it/it/opendata/Mappa_delle_stazioni_ferrovie_in_Italia)



# Train stations data / 2

Municipality	Province	Region	Name	Insertion Date	Insertion Date & Hour	OpenStreetMap ID	Longitude	Latitude
ALTRO	ALTRO	ALTRO	Albonago	2014	2014-11-11T16:22:37Z	984003073	8.9750146	46.0085186
ALTRO	ALTRO	ALTRO	Camedo	2014	2014-05-04T21:49:40Z	2836876118	8.610916	46.1547927
Alpignano	TORINO	Piemonte	Alpignano	2013	2013-05-04T21:56:45Z	1588214497	7.5238245	45.0912816
Avigliana	TORINO	Piemonte	Avigliana	2015	2015-05-11T21:25:04Z	1598967748	7.4009405	45.0853893
Balangero	TORINO	Piemonte	Balangero	2011	2011-08-14T22:12:24Z	822596562	7.5194836	45.265646
Bardonecchia	TORINO	Piemonte	Bardonecchia	2015	2015-05-11T21:25:04Z	2043486315	6.7097395	45.0764633
Borgaro Torinese	TORINO	Piemonte	Borgaro Torinese	2015	2015-05-11T21:25:02Z	3511147159	7.6487388	45.1532959
Borgofranco d'Ivrea	TORINO	Piemonte	Borgofranco d'Ivrea	2012	2012-11-20T09:13:02Z	1640450736	7.8546427	45.5147371
Borgone Susa	TORINO	Piemonte	Borgone di Susa	2015	2015-05-11T21:25:05Z	1598967749	7.2360913	45.1233785
Bosconero	TORINO	Piemonte	Bosconero	2011	2011-08-12T13:48:27Z	1262195587	7.7618626	45.2675134
Brandizzo	TORINO	Piemonte	Brandizzo	2014	2014-02-06T10:42:06Z	2043486321	7.8411549	45.1787318
Bruzolo	TORINO	Piemonte	Bruzolo	2015	2015-05-11T21:25:05Z	1598967750	7.2010571	45.1312692
Bussoleno	TORINO	Piemonte	Bussoleno	2015	2015-05-11T21:25:05Z	1446542028	7.1448124	45.1396163
Caluso	TORINO	Piemonte	Caluso	2012	2012-11-20T09:12:59Z	1003779384	7.9025604	45.3031447
Caluso	TORINO	Piemonte	Rodallo	2012	2012-11-20T09:12:59Z	1640251816	7.8780343	45.2787068

- 2,999 train stations
- Information already usable in tabular format
- Use station names to build and scrape urls

# Train stations data / 2

Municipality	Province	Region	Name	Insertion Date	Insertion Date & Hour	OpenStreetMap ID	Longitude	Latitude
ALTRO	ALTRO	ALTRO	Albonago	2014	2014-11-11T16:22:37Z	984003073	8.9750146	46.0085186
ALTRO	ALTRO	ALTRO	Camedo	2014	2014-05-04T21:49:40Z	2836876118	8.610916	46.1547927
Alpignano	TORINO	Piemonte	Alpignano	2013	2013-05-04T21:56:45Z	1588214497	7.5238245	45.0912816
Avigliana	TORINO	Piemonte	Avigliana	2015	2015-05-11T21:25:04Z	1598967748	7.4009405	45.0853893
Balangero	TORINO	Piemonte	Balangero	2011	2011-08-14T22:12:24Z	822596562	7.5194836	45.265646
Bardonecchia	TORINO	Piemonte	Bardonecchia	2015	2015-05-11T21:25:04Z	2043486315	6.7097395	45.0764633
Borgaro Torinese	TORINO	Piemonte	Borgaro Torinese	2015	2015-05-11T21:25:02Z	3511147159	7.6487388	45.1532959
Borgofranco d'Ivrea	TORINO	Piemonte	Borgofranco d'Ivrea	2012	2012-11-20T09:13:02Z	1640450736	7.8546427	45.5147371
Borgone Susa	TORINO	Piemonte	Borgone di Susa	2015	2015-05-11T21:25:05Z	1598967749	7.2360913	45.1233785
Bosconero	TORINO	Piemonte	Bosconero	2011	2011-08-12T13:48:27Z	1262195587	7.7618626	45.2675134
Brandizzo	TORINO	Piemonte	Brandizzo	2014	2014-02-06T10:42:06Z	2043486321	7.8411549	45.1787318
Bruzolo	TORINO	Piemonte	Bruzolo	2015	2015-05-11T21:25:05Z	1598967750	7.2010571	45.1312692
Bussoleno	TORINO	Piemonte	Bussoleno	2015	2015-05-11T21:25:05Z	1446542028	7.1448124	45.1396163
Caluso	TORINO	Piemonte	Caluso	2012	2012-11-20T09:12:59Z	1003779384	7.9025604	45.3031447
Caluso	TORINO	Piemonte	Rodallo	2012	2012-11-20T09:12:59Z	1640251816	7.8780343	45.2787068

- 2,999 train stations
- Information already usable in tabular format
- Use station names to build and scrape urls



◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ 🔍 ↺

# Wikipydia: scrape html pages

```
1 # Import the required libraries
2 import requests
3 from bs4 import BeautifulSoup # new
4 # specify the url to be scraped: Stazione di Tavazzano
5 wikiurl = 'https://it.wikipedia.org/wiki/Stazione_di_Tavazzano'
6 # GET method: scrape the url
7 r = requests.get(wikiurl)
8 html_page = r.text
9 # read the page with BeautifulSoup
10 soup = BeautifulSoup(r.text, 'html.parser')
11 # scrape the title name, and the other info
12 station_name = soup.find('h1', {'id' : 'firstHeading'}).text # station name
13 latitude = soup.find('span', {'class' : 'latitude'}).text # latitude
14 longitude = soup.find('span', {'class' : 'longitude'}).text # longitude
15 start_date = soup.find('th', string='Attivazione').findNext('a').text # sta
16 lines = soup.find('th', string='Linee').findNext('a').text # train line
17 status = soup.find('th', string='Stato attuale').findNext('td').text # acti
18 print([station_name,latitude,longitude,start_date,lines,status])
```

# Write down in readable formats - CSV

Scrape [www.congress.gov/](http://www.congress.gov/) with `trumpy.py`

```
import csv
import requests
from bs4 import BeautifulSoup
from os.path import dirname, abspath

### CONGRESS 112 ###
basedir = dirname(dirname(abspath(__file__))) + '/congress'
o112 = basedir + '/output.csv'
hrurl = 'https://www.congress.gov/bill/112th-congress/house-bill/6729/cosponsors?r=1'
hrpage = requests.get(hrurl)
data = hrpage.text
soup = BeautifulSoup(data, 'lxml') #here we make BS read the html page
with open(o112, 'w') as f: #open the outputfile in .csv
    fwriter = csv.writer(f, delimiter = ';') #declare the writing object
    fwriter.writerow(['Congress', 'BillType', 'BillNumber', 'cosponsorhref']) #write headers
    for l in soup.findAll('table', { 'class' : 'item_table'}): #identify cosponsors' table
        for a in l.findAll('a', href = True, text = True, target = '_blank'): #cosponsors' entries
            cosponsorhref = a['href']
            print([112, 'HR', 6729, cosponsorhref])
            fwriter.writerow([112, 'HR', 6729, cosponsorhref])
f.close()
```

Run

# Wikipydia: automate the analysis of rail stations

```
1 # Import the required libraries
2 import requests
3 from bs4 import BeautifulSoup
4 import os
5 import csv
6 # define directories
7 basedir = '/home/gabriele/Desktop/Dropbox (Francesco Decarolis)/Lessons/Tor
8 stationdir = os.path.join(basedir, 'Mappa-delle-stazioni-ferroviarie-in-Ita
9 outputdir = os.path.join(basedir, 'railstation_data.csv')
10 wrong_counter = 0
11 with open(stationdir, 'r+', encoding='latin1') as f: # without the encoding
12     reader = csv.reader(f, delimiter=';') # in this file's case, the delimi
13     with open(outputdir, 'w+') as g:
14         writer = csv.writer(g)
15         writer.writerow(['Comune', 'Provincia', 'Regione', 'Nome', 'Longitudine
16         for i, line in enumerate(reader): # going through all stations in t
17             # DOWNLOAD ROUTINE #
18     g.close()
19 f.close()
```

# Wikipydia: automate the analysis of rail stations

```
1 for i, line in enumerate(reader): # going through all stations in the list
2     if i <= 17:
3         continue # skip the first iterations: title + foreign stati
4     station_name = line[3] # line is the (iterable) .csv row: the 4
5     url_name = station_name.replace(' - ', '-').replace(' ', '_') #
6     wikiurl = 'https://it.wikipedia.org/wiki/Stazione_di_' + url_na
7     print(wikiurl)
8     # GET method: scrape the url
9     r = requests.get(wikiurl)
10    # read the page with BeautifulSoup
11    soup = BeautifulSoup(r.text, 'html.parser')
```

# Wikipydia: automate the analysis of rail stations

```
1 try:
2     status = soup.find('th', string='Stato attuale').findNext('td').text
3 except AttributeError:
4     status = '-'
5 try:
6     start_date = soup.find('th', string='Attivazione').findNext('a').text
7 except AttributeError:
8     start_date = '.'
9 try:
10    lines = soup.find('th', string='Linee').findNext('a').text
11 except AttributeError:
12    lines = '-'
13 if lines == '-' and start_date == '.' and status == '-': # the link is broken
14     lines = 'CHECK'
15     start_date = 'CHECK'
16     status = 'CHECK'
17     wron_counter += 1
18 printline = line[0:4]+line[7:9]+[status]+[start_date]+[lines]+[wikiurl] # v
19 writer.writerow(printline)
```

# Scraped Data

Municipality	Province	Region	Name	Longitude	Latitude	Status	Opening Date	Lines
Airasca	TORINO	Piemonte	Airasca	7.4832665	44.9280462	in uso	1854	Ferrovia Torino-Torre Pellice
Alpignano	TORINO	Piemonte	Alpignano	7.5238245	45.0912816	in uso	0	Frejus
Avigliana	TORINO	Piemonte	Avigliana	7.4009405	45.0853893	in uso	0	Frejus
Balangero	TORINO	Piemonte	Balangero	7.5194836	45.265646	In uso	1876	Ferrovia Torino-Ceres
Bardonecchia	TORINO	Piemonte	Bardonecchia	6.7097395	45.0764633	in uso	1871	Ferrovia del Frejus
Borgaro Torinese	TORINO	Piemonte	Borgaro Torinese	7.6487388	45.1532959	In uso	1871	Ferrovia Torino-Ceres
Borgofranco d'Ivrea	TORINO	Piemonte	Borgofranco d'Ivrea	7.8546427	45.5147371	CHECK	CHECK	CHECK
Borgone Susa	TORINO	Piemonte	Borgone di Susa	7.2360913	45.1233785	CHECK	CHECK	CHECK
Bosconero	TORINO	Piemonte	Bosconero	7.7618626	45.2675134	in uso	0	Canavesana
Brandizzo	TORINO	Piemonte	Brandizzo	7.8411549	45.1787318	In uso	1856	Torino-Milano
Bruzolo	TORINO	Piemonte	Bruzolo	7.2010571	45.1312692	in uso	0	Frejus
Bussoleno	TORINO	Piemonte	Bussoleno	7.1448124	45.1396163	in uso	1854	Torino-Modane
Caluso	TORINO	Piemonte	Caluso	7.9025604	45.3031447	in uso	0	Chivasso-Aosta
Caluso	TORINO	Piemonte	Rodallo	7.8780343	45.2787068	in uso	0	Chivasso-Aosta
Cambiano	TORINO	Piemonte	Cambiano-Santena	7.7731065	44.9664372	in uso	1849	Torino-Genova
Candia Canavese	TORINO	Piemonte	Candia Canavese	7.891791	45.3300205	in uso	0	Chivasso-Aosta
Candiolo	TORINO	Piemonte	Candiolo	7.5990181	44.9611065	in uso	1854	Ferrovia Torino-Torre Pellice
Carmagnola	TORINO	Piemonte	Carmagnola	7.7278083	44.8470455	in uso	1853	Torino-Savona
Caselle Torinese	TORINO	Piemonte	Caselle Aeroporto	7.6416415	45.1918384	In uso	2001	Ferrovia Torino-Ceres

- Works pretty well!
- Still needs refinements: station names for URLs, multiple lines, etc.
- Is the station list comprehensive?

# Requests - POST method

Scrape [www.semrush.com](http://www.semrush.com) with `pygliaSEM.py`

```
import requests
# list credentials and generate the login dict
username = "myloginmail@notmail.com"
password = "verysecretpassword"
login_data = {
    'email': username,
    'password': password
}
# define the base URL --> the one in which we need to login
base_url = 'https://www.semrush.com/json_users/login'
# start a request session and POST the login data
s = requests.session() #declare a SESSION of requests
r = s.post(base_url, data = login_data) #this way the session is logged in!
# proceed in parsing the .csv file
r = __csvParse__(domain,s,database,timestamp, 'domain_adwords')
```



# Requests - GET method with parameters

## Scrape www.semrush.com with pygliaSEM.py

```
def __csvParse__(domain,s,database,timestamp,domain_type):
    domain_url = 'https://www.semrush.com/info/' + domain + 'ad_domain'
    try:
        # initialize parameter dict and launch GET request
        params = {
            'action' : 'report_tpye',
            'database' : 'domain_name',
            'rnd_m' : 'rnd_number',
            'key' : 'some_key',
            'domain' : domain,
            'type' : 'domain_whatever',
            'display_filter' : '',
            'export_hash' : 'hash_key',
            'display_sort' : 'nq_desc',
            'export_columns' : 'column_list',
            'export_decode' : '',
            'export_escape' : '',
            'exchange_rate' : '',
            'currency' : 'usd',
            'export' : 'csv'
        }
        # download the .csv
        r = s.get('https://' + timestamp + database + '.backend.semrush.com/', params = params)
    except:
        r = ''
    return r
```

# Webdriver - when URLs are not enough

## Scrape www.coeweb.istat.it with pystat.py

```
from selenium import webdriver
from selenium.webdriver.common.keys import Keys
from selenium.webdriver.support.ui import WebDriverWait
import time
from os.path import dirname, abspath

##### ROUTINE #####
search_url = 'https://www.coeweb.istat.it/predefinite/tutto_paese_merce.asp?livello=L2&riga=MERCE&territorio=S&'
basedir = dirname(dirname(abspath(__file__))) + "/coeweb"
downloadedname = basedir + '/tabella.xls'

# FF settings
profile = webdriver.FirefoxProfile()
profile.set_preference('browser.download.folderList', basedir) # custom location
profile.set_preference('browser.download.dir', basedir)
profile.set_preference('browser.helperApps.neverAsk.saveToDisk', 'text/csv')
browser = webdriver.Firefox(profile)
browser.implicitly_wait(6) # wait some seconds for uploading the page
browser.get(search_url) # go to the page to be scraped
browser.find_element_by_name("XLS").click() # choose Excel format
browser.find_element_by_name("MESE").send_keys('Ultimo Disponibile') # Month
browser.find_element_by_id("radio5").click() # Cumulative
browser.find_element_by_name("B2").click() # download
time.sleep(10)
browser.close()
```

Run

# Countless Applications!

- Istat - Coeweb [Code](#)
- Consiglio di Stato - Pareri [Code](#)
- Corte dei Conti - Sentenze [Code](#)
- Repubblica - Articles [Code](#)
- Corriere della Sera - Articles [Code](#)
- American Congress - all bills / sponsors / cosponsors [Code](#)
- Californian Schools - [repolist](#)

# Part 1: webscraping with requests

The goal of the python code will be to automatically download and save to local disk all .xls files available at <http://www.cde.ca.gov/ds/sp/ai/>. These are clustered in three groups:

- SAT test results 1998-2016
- ACT test results 1998-2016
- AP test results 1998-2016

each consisting of 18 individual .xls files. The scraper should parse the webpage (with `requests.get()`), identify the links of the files (with `BeautifulSoup()`), and save the contents of each individual URL in a distinct file (with `requests.get()` and `.write()`).

# Tips

- Pay attention to the embedded URLs: from 2015-2016 to 2013-2014 they are complete, whereas from 2012-2013 on they only store the final “branches”, which need the base url to properly work!
- You may find useful to use the names provided in the URL to save the relative file on your local disk: the `basename()` function from `os.path` will help:

```
from os.path import basename
```

# Solution 2 - 1

```
# define the lists to be compared
a = [1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89]
b = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13]
# for loop, if, and, or solution
import sys
def ls_overlap(A,B):
    if not isinstance(A,list) or not isinstance(B, list): # check that both inputs are lists
        sys.exit("both inputs must be lists!")
    res = []
    for e in set(A): # for every unique element in A
        if e in B and B not in res: # if A element is in B and not already in res
            res.append(e)
    return(res)
res = ls_overlap(a,b)
print(res)
# use "set()" and "intersection()"
res = list(set(a).intersection(set(b)))
print(res)
```

back

## Solution 2 - 2

*# define the list*

```
ls_a = [1, 4, 9, 16, 25, 36, 49, 64, 81, 100]
```

*# one-line solution: list comprehension, for, in and modulus*

```
res = [x for x in ls_a if x % 2 == 0]
```

```
print(res)
```

*# for loop, in, if and append*

```
res = []
```

```
for a in ls_a:
```

```
    if a % 2 == 0:
```

```
        res.append(a)
```

```
print(res)
```

back