TECHNOLOGIES FOR INFORMATION SYSTEMS ROJECT

PROF. LETIZIA TANCA

# *WEB SCRAPING OF ANSA.IT*

Alessandro Artoni - 899343

ACADEMIC YEAR 2017/2018



**POLITECNICO**

MILANO 1863

# Introduction

The purpose of this project is to collect economical articles from the website [ansa.it](ansa.it) using an automated process implemented with a script in python.

The script makes a request for a certain page, downloads it, takes the information that are needed and loops - making another request - until it reached the last article.

# The Process

This section describes the steps followed to reach the objective of the project.

## Landing page analysis of Ansa

Firstly, the landing page of Ansa was analysed. We noticed that there is a form where an user can make a post to the database, and be redirect to a "result-page".

At this point, we needed to craft a proper post request, based on the attributes that the form needs.
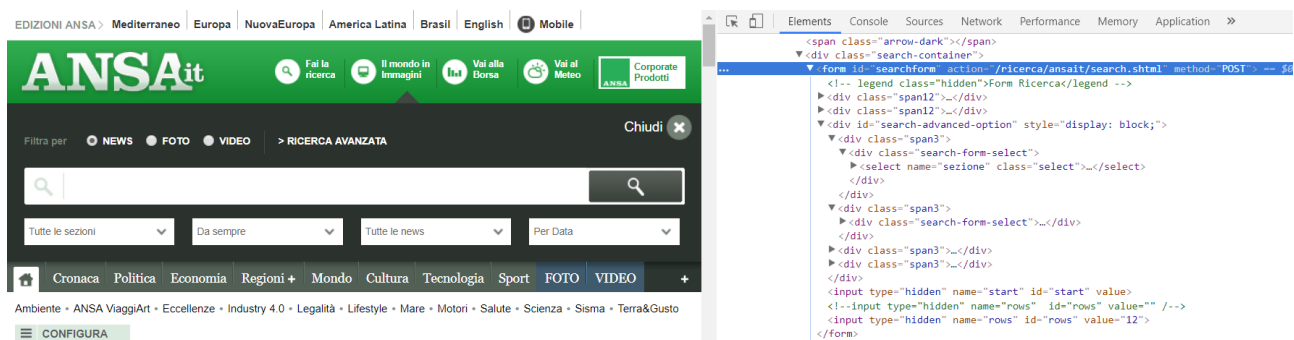


*Figure 1 - Analysis of ANSA landing page*

From Fig. 1, we can see the fields needed, and from Fig. 2, the corresponding post request.

The first span12, is a container for the type of the news (news, foto, video). If it is left empty, it selects news.

The second span12, is a container for the search request. In the program there's "company" which is a variable that contains the name of the company from wich we wish to retrieve information.

"Sezione" means section. Since we were interested in "Economia", we selected '63a85942-dedb-4b31-a3bf-a06f721c67e6' which should be encrypted but it's actually written clear-text on the source html. Other fields display how do we want to sort dates (in our case from the most recent to the least recent), and the number of rows we wanted to display.

```
104     # section '...' is actually 'Economia', other sections were encrypted. check www.ansa.com for other sections
105     post_fields = {'tiponotizia': '',
106                     'any': company,
107                     'sezione': '63a85942-dedb-4b31-a3bf-a06f721c67e6',
108                     'periodo': '',
109                     'genere': '',
110                     'sort': 'data:desc',
111                     'start': '',
112                     'rows': '12'}
```

*Figure 2 - Post request fields*

## Research page

Then, we analysed the next research page, and from that we started to retrieve some information.



*Figure 3 – Research Page*

In *red* we see the number of results, in *pink* there's the date and the time, in *orange* there's the title and in *blue* there's the subtitle.

From this page on, we retrieve the links of each article in the page. Then we make a simple get request to the article and we reach the article page.

## Article page



In the article page, we can double check that there is still the same title, a new subtitle and the body of the article itself.

Finally, it is important to note that the "number of results" is important because from that we know how many times we need to iterate in order to retrieve all the articles talking about a particular company.

## Query

Once we retrieved all these information, we clean all data we collected from weird html characters, in order to make them using several times the function `replace( html char, ascii – char)`.

After the clean, we can finally insert them into our database.

```
196     try:
197         with connection.cursor() as cursor:
198             query = "INSERT INTO articles_ansa (date, newspaper, section, title, summary,  " \
199                     "body, company, link_page) VALUES (%s, %s, %s, %s, %s, %s, %s, %s)"
200             cursor.execute(query, [date, "Ansa", category, title, abstract, body, company, link])
201             connection.commit()
202
203     except Exception, e:
204         print("Can't insert: logging in the file " + str(e))
```

*Figure 4 – Final query*

## Final remarks on the project

- There are no explicit authors on ANSA due to the fact that it is an agency. From an article page you will always see "*Riproduzione riservata © Copyright ANSA*"
- Some articles, marked "professional" are unavailable. You need to get your credential from ANSA as a journalist in order to read them.
- At the end of the project, we noticed that we saved a "string" date instead of a datetime type. Since it could be useful make some queries on the most recent news,  we provided python script to fix the problem. The python script though, needs a user to add a column on the database. The process could be fully automate but we preferred to avoid the script to drop/create columns inside tables.

## Conclusions

Thanks to this script, we provided an automated way to collect and store articles from ansa.it. At the end, more then 38 800 articles were registered, with all their corresponding information.