EUROPYTHON 2011

FLORENCE, JUNE 20-26

# Scraping Techniques to Extract Advertisements From Web Pages
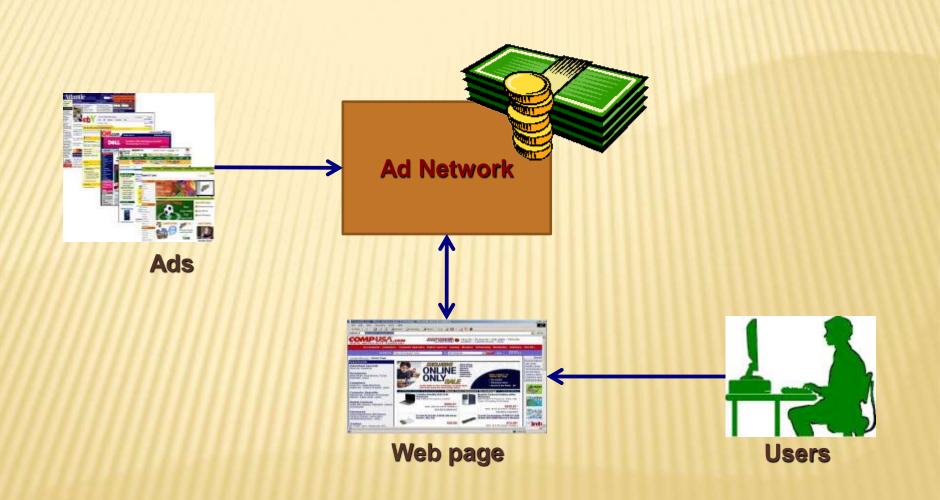
**Mirko Urru**
mirko.urru@hotmail.it

**Stefano Cotta Ramusino**
whitone@gmail.com

# OutLine

- Scraping techniques applied to contextual advertising
- Background
- What is scraping?
- The form of scraping
- The most  famous scraping techniques
- Application of scraping techniques
  to contextual advertising
- Conclusions

Ads

Ad Network

Web page

Users

# Sponsored Search

# Contextual Advertising

# Sponsored Search

# Sponsored Search



Web Site, Design

Web Site Design

**Origin:**

scraping data from mainframes
from green texts on black screens
to new data structures or API

**Nowadays:**

forcing data from old websites in
something new (web ≥ 2.0)

**Web scraping is the process of automatically collecting Web information**

**Beautiful Soup**

**mechanize**

**lxml**

**html5lib**

**scrapemark**

**pyquery**

**scrapy**

**QtWebKit**
**HtmlUnit**
**Selenium**
**Crowbar**
**Chrome Remote Shell**
**PyAuto**
**FireWatir**
**Spidermonkey**

# Beautiful Soup

www.crummy.com/software/BeautifulSoup

```python
soup = BeautifulSoup(webpage)

form = soup.find(type="password").findPrevious("form")

tag_input = form.findAll('input')

for tag in tag_input:
    name = tag['name']
    value = tag['value']

    inputs[name] = value
```

# Mechanize

wwwsearch.sourceforge.net/mechanize

```python
from mechanize import Browser

br = Browser()
br.open(uri)
assert br.viewing_html()
br.select_form(name="login")

br["username"] = "utente"
br["password"] = "segreto"

br.submit()
```

Inlink Extractor

{inlink}

Scraping
Module

{ads}

**Given a generic page the module extracts the p inlink Each inlink is displayed with the title and url**

titolo

url

```
href="http://www.crastulo.it/appuntamenti/1624_summer+beach+fest+2011.html"><img
src="/media/16/90050812922249/summer-beach.jpg" alt="SUMMER BEACH" height="177" width="615" /></a></div>
        <a class="thumb-btn-top-event" href="javascript:void(0);"><img src="/media/16/94942472976816/summer-
beach.jpg" alt="SUMMER BEACH" height="41" width="93" /></a>
</div>
<div class="thumb-top-event">
        <div class="image-top-event"><a
href="http://www.crastulo.it/appuntamenti/1700_l%27aperivizi+dei+sette+vizi.html"><img
src="/media/14/77734406237187/aperivizi.jpg" alt="SETTE VIZI" height="177" width="615" /></a></div>
        <a class="thumb-btn-top-event" href="javascript:void(0);"><img src="/media/15/45182478745611/sette-vizi.jpg"
alt="SETTE VIZI" height="41" width="93" /></a>
</div>
<div class="thumb-top-event">
        <div class="image-top-event"><a
href="http://www.crastulo.it/appuntamenti/1642_il+25+giugno+inaugura+il+blufan.html"><img
src="/media/15/87553322647925/blufan.jpg" alt="BLUFAN" height="177" width="615" /></a></div>
        <a class="thumb-btn-top-event" href="javascript:void(0);"><img src="/media/15/14759601997811/blufan.jpg"
alt="BLUFAN" height="41" width="93" /></a>
</div>
<div class="thumb-top-event">
        <div class="image-top-event"><a
href="http://www.crastulo.it/appuntamenti/1480_gli+eventi+dell%27estate+2011.html"><img
src="/media/15/46870991311486/jovanotti.jpg" alt="concerti" height="177" width="615" /></a></div>
        <a class="thumb-btn-top-event" href="javascript:void(0);"><img src="/media/15/42598090854818/jovanotti.jpg"
alt="concerti" height="41" width="93" /></a>
</div>
```

# I Scraping the Source Code

```python
tobeparsed = mechanize.urlopen(url)
body = BeautifulSoup.BeautifulSoup(tobeparsed)
body = body.prettify()
body = body[:body.find('</div>')]
link = body[body.find('<a')+len('</a>'):body.find('<br />')]
control = body[body.rfind('-->')+len('--
>'):body.find('</strong>')]
link = link.replace('&#039;',"'")
link = link.split()
control = control.split()
```

**The extracted ads are then randomly selected and displayed in the target Web Page**

# **Extracting inlink from : www.crastulo.it**

# Here's the four randomly selected ads

- This project was aimed at suggesting suitable ads to a given Web page
- To this end I devised a system written in Python that:
  - extracts a set of inlinks of a given Web page
  - randomly selects four ads previously extracted by scraping

- To apply  scraping techniques also for dynamic advertising
- To suggest ads according to users interests

## **Thanks to all**

Contact us for details and questions on scraping in Python:

mirko.urru@hotmail.it

whitone@gmail.com

Contact Eloisa Vargiu for details and questions on contextual advertising :

vargiu@diee.unica.it

**Thanks to**


Università di Cagliari