

Python for applied practitioners

F. Curci, F. Masera, T. Rodríguez Martínez

February 24th, 2017

Sketch presentation

How Python can be useful to applied practitioners

- Fede C
 - **What Python is and scraping**
- Fede M
 - GIS (geography and maps) application
- Tomás
 - Machine learning tools

What Python is

- Python is a widely used programming language for general-purpose programming. It has a syntax which allows programmers to express concepts in fewer lines of code than possible in languages such as C++ or Java
- Free
- Users contribution: libraries
- Many platforms use Python as communicating way (becoming universal language)
- Very powerful in dealing with strings

Good way to learn Python

- Coursera course: *Python for everybody specialization. University of Michigan*
- Book: *Python for informatics*

How to obtain it and language

- *Download*
- Languages currently used: Python 2 and Python 3
- *Anaconda*

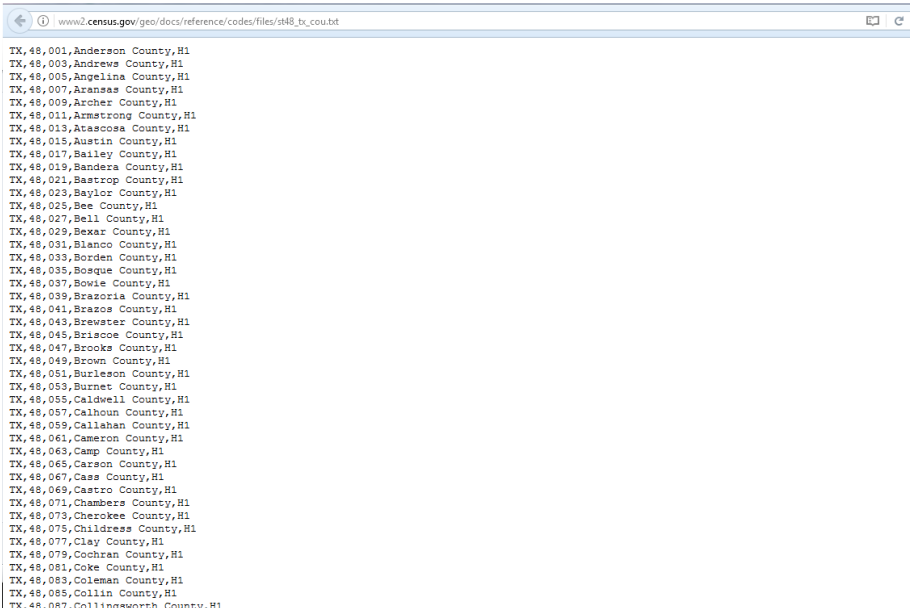
Two methods

- Notepad+command prompt
 - Need to save file as .py
 - Type file.py in command prompt
- Compiler: *Pycharm, Atom, Spyder*

Python allows you to write a web browser

- Library to open resources by URL: urllib
- Example: we want to Python to go to web to obtain FIPS code of Texas counties
 - *Counties*

Web browsing



Web browsing

The screenshot displays the PyCharm IDE interface. At the top, the title bar indicates the file path: `urlopen_county.py - [C:\Users\fcuirci\AppData\Local\Temp\urlopen_county.py] - C:\Users\fcuirci\Documents\UC3MVAppliedReading\Python\Codes\urlopen_county.py - PyCharm Community Edition 2016.3.2`. The menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, and Help. The Project tool window on the left shows the project structure with `urlopen_county.py` and `External Libraries`. The main editor window displays the following Python code:

```
1 import urllib
2 counties=urllib.urlopen('http://www2.census.gov/geo/docs/reference/codes/files/st40_tx_cou.txt')
3 f = line in counties:
4     print line.strip()
```

A yellow notification bar states: "No Python interpreter configured for the project". The Run tool window at the bottom shows the output of the script, listing Texas counties and their FIPS codes:

```
TX,48,419,Shelby County,H1
TX,48,421,Sherman County,H1
TX,48,423,Smith County,H1
TX,48,425,Somervell County,H1
TX,48,427,Starr County,H1
TX,48,429,Stephens County,H1
TX,48,431,Sterling County,H1
TX,48,433,Stonewall County,H1
TX,48,435,Sutton County,H1
TX,48,437,Swisher County,H1
TX,48,439,Tarrant County,H1
TX,48,441,Taylor County,H1
TX,48,443,Terrell County,H1
TX,48,445,Terry County,H1
TX,48,447,Throckmorton County,H1
TX,48,449,Titus County,H1
TX,48,451,Tom Green County,H1
TX,48,453,Travis County,H1
TX,48,455,Trinity County,H1
TX,48,457,Tyler County,H1
TX,48,459,Upshur County,H1
TX,48,461,Upton County,H1
TX,48,463,Uvalde County,H1
TX,48,465,Val Verde County,H1
TX,48,467,Van Zandt County,H1
TX,48,469,Victoria County,H1
TX,48,471,Walker County,H1
TX,48,473,Waller County,H1
TX,48,475,Ward County,H1
TX,48,477,Washington County,H1
TX,48,479,Webb County,H1
TX,48,481,Wharton County,H1
```


Web browsing

The screenshot shows the PyCharm IDE interface. At the top, the title bar indicates the file path: `urlopen_county.py - [C:\Users\fcuirci\AppData\Local\Temp\urlopen_county.py] - C:\Users\fcuirci\Documents\UC3MVAppliedReading\Python\Codes\urlopen_county.py - PyCharm Community Edition 2016.3.2`. The menu bar includes File, Edit, View, Navigate, Code, Refactor, Run, Tools, VCS, Window, and Help. The toolbar shows icons for Project, Run, and other development actions. The left sidebar displays the Project view with `urlopen_county.py` and External Libraries. The main editor window shows the code for `urlopen_county.py` with a yellow highlight around the `import urllib` line. A message box states "No Python interpreter configured for the project". The Run console at the bottom shows the output of the script, listing Texas counties and their FIPS codes, with `TX, 48, 453, Travis County, H1` highlighted.

```
urlopen_county.py - [C:\Users\fcuirci\AppData\Local\Temp\urlopen_county.py] - C:\Users\fcuirci\Documents\UC3MVAppliedReading\Python\Codes\urlopen_county.py - PyCharm Community Edition 2016.3.2
File Edit View Navigate Code Refactor Run Tools VCS Window Help
urlopen_county.py
Project urlopen_county.py
urlopen_county.py
No Python interpreter configured for the project
1 import urllib
2 counties=urllib.urlopen('http://www2.census.gov/geo/docs/reference/codes/files/st48_tx_cou.txt')
3 for line in counties:
4     print line.strip()
Run urlopen_county
TX, 48, 419, Shelby County, H1
TX, 48, 421, Sherman County, H1
TX, 48, 423, Smith County, H1
TX, 48, 425, Somervell County, H1
TX, 48, 427, Starr County, H1
TX, 48, 429, Stephens County, H1
TX, 48, 431, Sterling County, H1
TX, 48, 433, Stonewall County, H1
TX, 48, 435, Sutton County, H1
TX, 48, 437, Swisher County, H1
TX, 48, 439, Tarrant County, H1
TX, 48, 441, Taylor County, H1
TX, 48, 443, Terrell County, H1
TX, 48, 445, Terry County, H1
TX, 48, 447, Throckmorton County, H1
TX, 48, 449, Titus County, H1
TX, 48, 451, Tom Green County, H1
TX, 48, 453, Travis County, H1
TX, 48, 455, Trinity County, H1
TX, 48, 457, Tyler County, H1
TX, 48, 459, Upshur County, H1
TX, 48, 461, Upton County, H1
TX, 48, 463, Uvalde County, H1
TX, 48, 465, Val Verde County, H1
TX, 48, 467, Van Zandt County, H1
TX, 48, 469, Victoria County, H1
TX, 48, 471, Walker County, H1
TX, 48, 473, Waller County, H1
TX, 48, 475, Ward County, H1
TX, 48, 477, Washington County, H1
TX, 48, 479, Webb County, H1
TX, 48, 481, Wharton County, H1
```

Ways to extract information from the web

API (Application programming interface)

- Structured way to access data from websites
- Usually limited number of queries
- Not all website provide it
- Twitter, Google maps, etc.
- Some Stata commands

Scraping

- Transformation of unstructured data (HTML format) on the web into structured data (database)
- Stata command: copy

Talk to web services: URLs that are designed explicitly to hand data back cleaned for your application

- Need to agree on wire protocol for python and java

Two wire formats to exchange data between applications

- XML
- JSON

Libraries

- `xml.etree.ElementTree`
- `json`

Google map API

- Possible to retrieve address, latitude and longitude of places
- *Documentation*
- Limits
 - 2,500 requests per day
 - 50 requests per second

Google map API

geojson_loop.py - [C:\Users\tcurch\AppData\Local\Temp\geojson_loop.py] - C:\Users\tcurch\Documents\UC3M\AppliedReading\Python\Codes\geojson_loop.py - PyCharm Community Edition 2016.3.2

File Edit View Navigate Code Refactor Run Tools VCS Window Help

geojson_loop.py

Project

geojson_loop.py

External Libraries

No Python interpreter configured for the project

```
1 import urllib
2 import json
3
4 universidades=["Universidad de Alcala","Universidad Autonoma de Madrid","Universidad Carlos III de Madrid",
5
6 serviceurl = 'http://maps.googleapis.com/maps/api/geocode/json?'
7 # serviceurl = 'http://python-data.dr-chuck.net/geojson?'
8
9 for universidad in universidades:
10     address = universidad
11     url = serviceurl + urllib.urlencode({'sensor':'false', 'address': address})
12     print 'Retrieving', url
13     uh = urllib.urlopen(url)
14     data = uh.read()
15     try: js = json.loads(str(data))
16     except: js = None
17     if 'status' not in js or js['status'] != 'OK':
18         print '==== Failure To Retrieve ====='
19         continue
20     lat = js["results"][0]["geometry"]["location"]["lat"]
21     lng = js["results"][0]["geometry"]["location"]["lng"]
22     print 'lat',lat,'lng',lng
23     location = js["results"][0]["formatted_address"]
24     print universidad +',' + location
```

Run geojson_loop

```
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+de+Alcala
lat 40.4824722 lng -3.3628674
Universidad de Alcala,Plaza de San Diego, s/n, 28801 Alcalá de Henares, Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Autonoma+de+Madrid
lat 40.5466983 lng -3.6943619
Universidad Autonoma de Madrid,Ciudad Universitaria de Cantoblanco, 28049 Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Carlos+III+de+Madrid
lat 40.316966 lng -3.7270795
Universidad Carlos III de Madrid,Calle Madrid, 126, 28903 Getafe, Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Complutense+de+Madrid
lat 40.4478246 lng -3.7285872
```

Google map API

geojson_loop.py - [C:\Users\tcurch\AppData\Local\Temp\geojson_loop.py] - C:\Users\tcurch\Documents\UC3M\AppliedReading\Python\Codes\geojson_loop.py - PyCharm Community Edition 2016.3.2

File Edit View Navigate Code Refactor Run Tools VCS Window Help

geojson_loop.py

Project

geojson_loop.py

External Libraries

No Python interpreter configured for the project

```
1 import urllib
2 import json
3
4 universidades=["Universidad de Alcala","Universidad Autonoma de Madrid","Universidad Carlos III de Madrid"]
5
6 serviceurl = 'http://maps.googleapis.com/maps/api/geocode/json?'
7 # serviceurl = 'http://python-data.dr-chuck.net/geojson?'
8
9 for universidad in universidades:
10     address = universidad
11     url = serviceurl + urllib.urlencode({'sensor':'false', 'address': address})
12     print 'Retrieving', url
13     uh = urllib.urlopen(url)
14     data = uh.read()
15     try: js = json.loads(str(data))
16     except: js = None
17     if 'status' not in js or js['status'] != 'OK':
18         print '==== Failure To Retrieve ====='
19         continue
20     lat = js["results"][0]["geometry"]["location"]["lat"]
21     lng = js["results"][0]["geometry"]["location"]["lng"]
22     print 'lat',lat,'lng',lng
23     location = js["results"][0]["formatted_address"]
24     print universidad +','+ location
```

Run geojson_loop

```
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+de+Alcala
lat 40.4824722 lng -3.3628674
Universidad de Alcala,Plaza de San Diego, s/n, 28801 Alcalá de Henares, Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Autonoma+de+Madrid
lat 40.5466983 lng -3.6943619
Universidad Autonoma de Madrid.Ciudad Universitaria de Cantoblanco. 28049 Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Carlos+III+de+Madrid
lat 40.316966 lng -3.7270795
Universidad Carlos III de Madrid,Calle Madrid, 126, 28903 Getafe, Madrid, Spain
Retrieving http://maps.googleapis.com/maps/api/geocode/json?sensor=false&address=Universidad+Complutense+de+Madrid
lat 40.4478246 lng -3.7285872
```

HTML

- Standard markup language for creating web pages and web applications. Messy

```
<html>
  <head>
    <title>This is the title</title>
  </head>
  <body>
    <a href ="www.wikipedia.es" > Wikipedia link</a>
  </body>
</html>
```

Web scraping

HTML

- Standard markup language for creating web pages and web applications. Messy

`<html>`

`<head>`

Start tag

`<title>This is the title</title>`

Element

`</head>`

End tag

`<body>`

`Wikipedia link`

Attribute name

Attribute value

Content

`</body>`

`</html>`

Web scraping

HTML

- Standard markup language for creating web pages and web applications. Messy

`<html>`

`<head>`

Start tag

`<title>This is the title</title>`

Element

`</head>`

End tag

`<body>`

`Wikipedia link`

Attribute name

Attribute value

Content

`</body>`

`</html>`

- Possible to retrieve all html with python
- Libraries: *BeautifulSoup*, *HTMLParser*

← ⓘ https://www.idealista.com/venta-viviendas/madrid-madrid/mapa ↻

✿ Pon tu an

64.692

[illegible]

 Dibujar tu propia zona



Y además 245 promos de obra nueva

Arganzuela (763)

Looking for house prices

https://www.idealista.com/venta-viviendas/madrid-madrid/mapa

Mapa de Madrid, ver las 27.054 viviendas »

Y además 245 promos de obra nueva

Áreas de Madrid

- Arganzuela (763)
- Barajas (399)
- Carabanchel (1.763)
- Centro (1.331)
- Chamartín (2.101)
- Chamberí (1.338)
- Ciudad Lineal (1.494)
- Fuencarral (1.494)
- Hortaleza (1.798)
- Latina (1.188)
- Moncloa (2.041)

Context menu for Arganzuela (763):

- Abrir enlace en una pestaña nueva
- Abrir enlace en una ventana nueva
- Abrir enlace en una nueva ventana privada
- Añadir este enlace a marcadores
- Guardar enlace como...
- Guardar enlace en Pocket
- Copiar la ruta del enlace
- Buscar "Arganzuela (763..." en Google
- Inspeccionar elemento

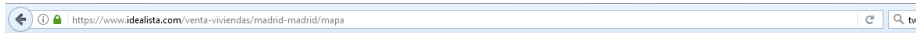
Developer Tools:

- Inspector
- Console
- Depurador
- () Editor de es...
- ⌘ Rendimiento
- ⌘ Memoria
- ⌘ Red

2 de 27

```
<div id="map-title-mobile" class="d-none"></div>
<div id="map-template" class="d-none"></div>
<div id="grey-map-wrapper">
  <map id="map-mapping" name="map-mapping">
    <area shape="poly" coords="126,284,132,284,132,294,136,300,142,305,155,308,169,308,175,...,334,162,331,153,328,146,323,140,317,127,317,126,312,126,284"
      href="/venta-viviendas/madrid/arganzuela/mapa" data="Arganzuela, 763 viviendas">
    <area shape="poly" coords="385,29,428,29,437,25,442,18,448,16,448,96,448,176,382,175,38...,103,333,97,337,92,346,86,356,75,365,56,383,38,385,31,385,29"
      href="/venta-viviendas/madrid/barajas/mapa" data="Barajas, 399 viviendas">
```

Looking for house prices



Mapa de Madrid, ver las 27.054 viviendas »

Y además 245 promos de obra nueva

Áreas de Madrid

[Arganzuela \(763\)](#)

[Barajas \(399\)](#)

[Carabanchel \(1.7\)](#)

[Centro \(1.331\)](#)

[Chamartín \(2.10\)](#)

[Chamberí \(1.338\)](#)

[Ciudad Lineal \(1\)](#)

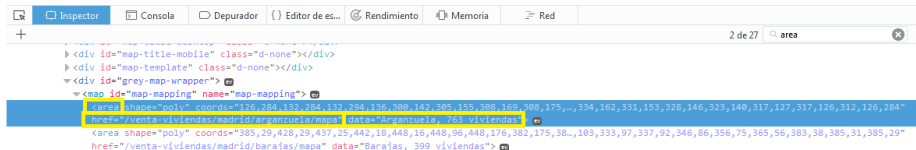
[Fuencarral \(1.494\)](#)

[Hortaleza \(1.798\)](#)






[Latina \(1.188\)](#)

[Moncloa \(2.041\)](#)

- Abrir enlace en una pestaña nueva
- Abrir enlace en una ventana nueva
- Abrir enlace en una nueva ventana privada
- Añadir este enlace a marcadores
- Guardar enlace como...
- Guardar enlace en Pocket
- Copiar la ruta del enlace
- Buscar "Arganzuela (763..." en Google
- Inspeccionar elemento



Looking for house prices

```
.....  
▶ <div id="map-title-mobile" class="d-none"></div>  
▶ <div id="map-template" class="d-none"></div>  
▼ <div id="grey-map-wrapper">   
  ▼ <map id="map-mapping" name="map-mapping">   
    <area shape="poly" coords="126,284,132,284,132,294,136,300,142,305,155,308,169,308,175,...,334,162,331,153,328,146,323,140,317,127,126,284" data-bbox="132 485 1000 500"/>  
    href="/venta-viviendas/madrid/arganzuela/mapa" data="Arganzuela, 763 viviendas"   
    <area shape="poly" coords="385,29,428,29,437,25,442,18,448,16,448,96,448,176,382,175,38...,103,333,97,337,92,346,86,356,75,365,56,385,29" data-bbox="132 525 1000 540"/>  
    href="/venta-viviendas/madrid/barajas/mapa" data="Barajas, 399 viviendas"   
    <area shape="poly" coords="23,433,39,414,44,397,52,379,43,368,41,365,46,364,57,370,62,3...419,112,419,97,425,86,426,73,426,52,433,23,433" data-bbox="132 570 1000 585"/>  
    href="/venta-viviendas/madrid/carabanchel/mapa" data="Carabanchel, 1.732 viviendas" 
```

Looking for house prices

File Edit View Navigate Code Refactor Run Tools VCS Window Help

2_idealista_barrio.py

Project
2_idealista_barrio.py

External Libraries

Structure

No Python interpreter configured for the project

```
6
7 import urllib
8 import time
9 from BeautifulSoup import *
10
11 url='https://www.idealista.com/venta-viviendas/madrid-madrid/mapa'
12 html = urllib.urlopen(url).read()
13 soup = BeautifulSoup(html)
14 tags = soup('area')
15 barrios=list()
16 href_barrios=list()
17 list_barrio_href=list()
18 for tag in tags:
19     colindancias=tag.get('h4', None)
20     barrio=tag.get('data', None)
21     firstpos=barrio.find(',')
22     secondpos=barrio.find('viviendas')
23     #Avoid to take neighbouring cities (colindancias)
24     if firstpos==len(barrio)-1:
25         continue
26     else:
27         barrios.append(barrio[0:firstpos])
28         href_barrio=tag.get('href', None)
29         href_barrios.append(href_barrio)
30         list_barrio_href.append([barrio[0:firstpos],barrio[firstpos+1:secondpos],href_barrio])
31
32 for barrio in list_barrio_href:
33     for minibarrio in barrio:
34         print minibarrio
```

Run 2_idealista_barrio

```
C:\Python27\python.exe C:/Users/fourci/Documents/UC3M/AppliedReading/Python/Codes/2_idealista_barrio.py
Arganzuela
763
/venta-viviendas/madrid/arganzuela/mapa
Barajas
399
/venta-viviendas/madrid/barajas/mapa
Carabanchel
1.732
/venta-viviendas/madrid/carabanchel/mapa
```

Looking for house prices

File Edit View Navigate Code Refactor Run Tools VCS Window Help

2_idealista_barrio.py

Project
2_idealista_barrio.py
External Libraries

Structure

No Python interpreter configured for the project

```
6
7 import urllib
8 import time
9 from BeautifulSoup import *
10
11 url='https://www.idealista.com/venta-viviendas/madrid-madrid/mapa'
12 html = urllib.urlopen(url).read()
13 soup = BeautifulSoup(html)
14 tags = soup('area')
15 barrios=list()
16 href_barrios=list()
17 list_barrio_href=list()
18 for tag in tags:
19     colindancias=tag.get('hd', None)
20     barrio=tag.get('data', None)
21     firstpos=barrio.find(',')
22     secondpos=barrio.find('viviendas')
23     #Avoid to take neighbouring cities (colindancias)
24     if firstpos==len(barrio)-1:
25         continue
26     else:
27         barrios.append(barrio[0:firstpos])
28         href_barrio=tag.get('href', None)
29         href_barrios.append(href_barrio)
30         list_barrio_href.append([barrio[0:firstpos],barrio[firstpos+1:secondpos],href_barrio])
31
32 for barrio in list_barrio_href:
33     for minibarrio in barrio:
34         print minibarrio
```

Run 2_idealista_barrio

C:\Python27\python.exe C:/Users/fourci/Documents/UC3M/AppliedReading/Python/Codes/2_idealista_barrio.py

```
Arganzuela
763
/venta-viviendas/madrid/arganzuela/mapa
barajas
399
/venta-viviendas/madrid/barajas/mapa
Carabanchel
1.732
/venta-viviendas/madrid/carabanchel/mapa
```

Looking for house prices

https://www.idealista.com/venta-viviendas/madrid/arganzuela/acacias/

Nuevos anuncios en tu email:

Guardar búsqueda

Viviendas

Precio

Mín Máx

Tamaño

Mín Máx

Tipo de vivienda

☐ Pisos

☐ Casas o chalets

☐ Casas rústicas

☐ Dúplex

☐ Áticos

Habitaciones

☐ 0 habitaciones (estudios)

Ordenar: **Relevancia** Baratos Recientes Más

Listado Mapa

Destacado

Piso en calle de Martín de Vargas, Acacias, Madrid

315.000 €

3 hab. 135 m² 4ª planta exterior con ascensor

"Inmobiliarias Encuentro vende piso todo exterior de 135 metros cuadrados, consta de 3 dormitorios amplios, 1 cuartos de baño comple...

912 664 063

Destacado

Piso en calle Majuelo, 3, Acacias, Madrid

429.000 € Garaje incluido

2 hab. 85 m² 2ª planta exterior con ascensor

De las mejores urbanizaciones de Pasillo Verde. Piso espectacular con calidades de primera. La cocina de diseño super equipada con tended...

916 350 423

Destacado

Dúplex en paseo de las acacias, 6, Acacias, Madrid

358.500 €

2 hab. 103 m² 5ª planta exterior con ascensor

Inspector **Consola** **Depurador** **Editor de es...** **Rendimiento** **Memoria** **Red**

Buscar en HTML

Reglas **Calculado**

Filtrar estilos

```
elemento {
  .row.price-row > span:last-child
    margin-right: 0;
}
.row.price-row .item-price {
  font-weight: bold;
  font-size: 1.375em;
  vertical-align: middle;
  display: inline-block;
  line-height: normal;
```

clearfix > div#main-content > div.items-container > article > div.item.item_contains_branding.item_fad... > div.clearfix > div.item-info-container > div.row.price-row.clearfix > span.item-price

Looking for house prices

https://www.idealista.com/venta-viviendas/madrid/arganzuela/acacias/

Nuevos anuncios en tu email:
Guardar búsqueda

Viviendas

Precio
Min Máx

Tamaño
Min Máx

Tipo de vivienda

- ☐ Pisos
- ☐ Casas o chalets
- ☐ Casas rústicas
- ☐ Dúplex
- ☐ Áticos

Habitaciones

- ☐ 0 habitaciones (estudios)

Ordenar: **Relevancia** Baratos Recientes Más

Destacado

Piso en calle de Martín de Vargas, Acacias, Madrid

315.000 €

3 hab. 135 m² 4ª planta exterior con ascensor

"Inmobiliarias Encuentro vende piso todo exterior de 135 metros cuadrados, consta de 3 dormitorios amplios, 1 cuartos de baño comple...

912 664 063

Destacado

Piso en calle Majuelo, 3, Acacias, Madrid

429.000 € Garaje incluido

2 hab. 85 m² 2ª planta exterior con ascensor

De las mejores urbanizaciones de Pasillo Verde. Piso espectacular con calidades de primera. La cocina de diseño super equipada con tended...

916 350 423

Destacado

Dúplex en paseo de las acacias, 6, Acacias, Madrid

358.500 €

2 hab. 103 m² 5ª planta exterior con ascensor

Inspector **Consola** **Depurador** **Editor de es...** **Rendimiento** **Memoria** **Red**

Buscar en HTML

Reglas **Calculado**

Filtrar estilos

elemento {

.row.price-row > span:last-child

margin-right: 0;

.row.price-row .item-price {

font-weight: bold;

font-size: 1.375em;

vertical-align: middle;

display: inline-block;

line-height: normal;

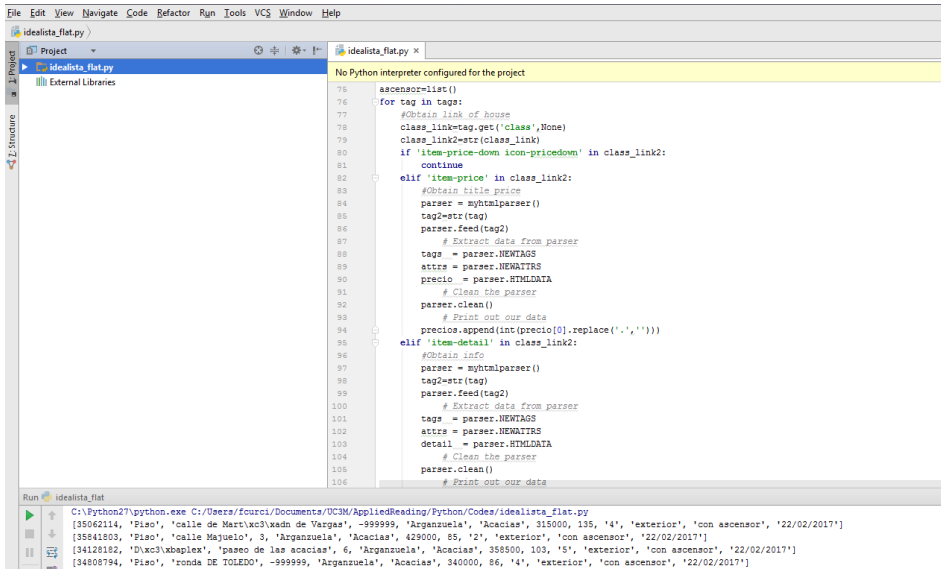
```
<div class="logo-branding"></div>
<div class="item-link" href="/inmueble/35862114/" title="Piso en calle de Martín de Vargas, Acacias, Madrid" data-xiti-click="listado:enlace">
  <div class="row price-row clearfix">
    <div>
      <span class="item-price">
        315.000
      </span>
    </div>
  </div>
</div>
```

clearfix > div#main-content > div.items-container > article > div.item.item_contains_branding.item_fad... > div.clearfix > div.item-info-container > div.row.price-row.clearfix > span.item-price

Looking for house prices

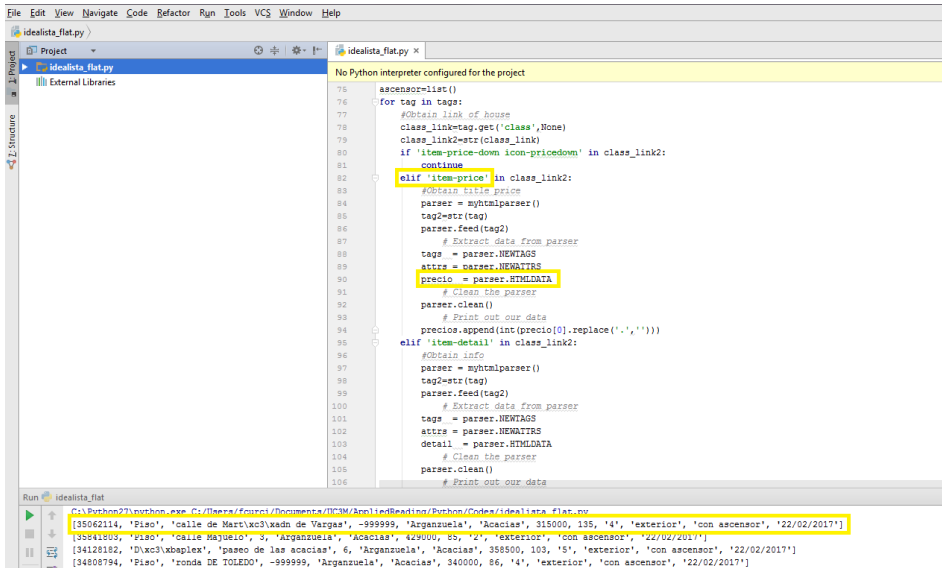
```
<div class="logo-branding"></div>
<a class="item-link" href="/inmueble/35062114/" title="Piso en calle de Martín de Vargas, Acacias, Madrid" data-xiti-click="li">
Piso en calle de Martín de Vargas, Acacias, Madrid</a>
<div class="row price-row clearfix">
  <div class="price">
    <span class="item-price">
      315.000
    </span>
  </div>
</div>
```

Looking for house prices



```
File Edit View Navigate Code Refactor Run Tools VCS Window Help
idealista_flat.py
Project
idealista_flat.py
External Libraries
No Python interpreter configured for the project
75 ascensor=list()
76 for tag in tags:
77     #Obtain link of house
78     class_link=tag.get('class',None)
79     class_link2=str(class_link)
80     if 'item-price-down icon-pricedown' in class_link2:
81         continue
82     elif 'item-price' in class_link2:
83         #Obtain title price
84         parser = myhtmlparser()
85         tag2=str(tag)
86         parser.feed(tag2)
87         # Extract data from parser
88         tags = parser.NEWTAGS
89         attrs = parser.NEWATTRS
90         precio = parser.HTMLDATA
91         # Clean the parser
92         parser.clean()
93         # Print out our data
94         precios.append(int(precio[0].replace('.', '')))
95     elif 'item-detail' in class_link2:
96         #Obtain info
97         parser = myhtmlparser()
98         tag2=str(tag)
99         parser.feed(tag2)
100         # Extract data from parser
101         tags = parser.NEWTAGS
102         attrs = parser.NEWATTRS
103         detail = parser.HTMLDATA
104         # Clean the parser
105         parser.clean()
106         # Print out our data
Run idealista_flat
C:\Python27\python.exe C:/Users/fourci/Documents/UC3M/AppliedReading/Python/Codes/idealista_flat.py
[35062114, 'Piso', 'calle de Mart\x03\xadn de Vargas', -999999, 'Arganzuela', 'Acacias', 315000, 135, '4', 'exterior', 'con ascensor', '22/02/2017']
[35841803, 'Piso', 'calle Majuelo', 3, 'Arganzuela', 'Acacias', 429000, 85, '2', 'exterior', 'con ascensor', '22/02/2017']
[34128182, 'D\x03\xbples', 'paseo de las acacias', 6, 'Arganzuela', 'Acacias', 358500, 103, '5', 'exterior', 'con ascensor', '22/02/2017']
[34808794, 'Piso', 'ronda DE TOLEDO', -999999, 'Arganzuela', 'Acacias', 340000, 86, '4', 'exterior', 'con ascensor', '22/02/2017']
```

Looking for house prices



```
File Edit View Navigate Code Refactor Run Tools VCS Window Help
idealista_flat.py
Project
idealista_flat.py
External Libraries
No Python interpreter configured for the project
75 ascensor=list()
76 for tag in tags:
77     #Obtain link of house
78     class_link=tag.get('class',None)
79     class_link2=str(class_link)
80     if 'item-price-down icon-pricedown' in class_link2:
81         continue
82     elif 'item-price' in class_link2:
83         #Obtain title price
84         parser = myhtmlparser()
85         tag2=str(tag)
86         parser.feed(tag2)
87         # Extract data from parser
88         tags = parser.NEWTAGS
89         attrs = parser.NEWATTRS
90         precio = parser.HTMLDATA
91         # Clean the parser
92         parser.clean()
93         # Print out our data
94         precios.append(int(precio[0].replace('.', '')))
95     elif 'item-detail' in class_link2:
96         #Obtain info
97         parser = myhtmlparser()
98         tag2=str(tag)
99         parser.feed(tag2)
100         # Extract data from parser
101         tags = parser.NEWTAGS
102         attrs = parser.NEWATTRS
103         detail = parser.HTMLDATA
104         # Clean the parser
105         parser.clean()
106         # Print out our data
C:\Python27\python.exe C:\Users\fourci\Documents\UIC3M\AppliedReading\Python\Codes\idealista_flat.py
[35062114, 'Piso', 'calle de Mart\x03\xadn de Vargas', -999999, 'Arganzuela', 'Acacias', 315000, 135, '4', 'exterior', 'con ascensor', '22/02/2017']
[35841803, 'Piso', 'calle Majuelo', 3, 'Arganzuela', 'Acacias', 429000, 85, '2', 'exterior', 'con ascensor', '22/02/2017']
[34128182, 'D\x03\xbples', 'paseo de las acacias', 6, 'Arganzuela', 'Acacias', 358500, 103, '5', 'exterior', 'con ascensor', '22/02/2017']
[34808794, 'Piso', 'ronda DE TOLEDO', -999999, 'Arganzuela', 'Acacias', 340000, 86, '4', 'exterior', 'con ascensor', '22/02/2017']
```

- Easier when PDF is not an image
- Possible to convert entire PDF to txt file
- Possible to dump the internal contents of a PDF file in pseudo-XML format
- Library: *PDFMiner*

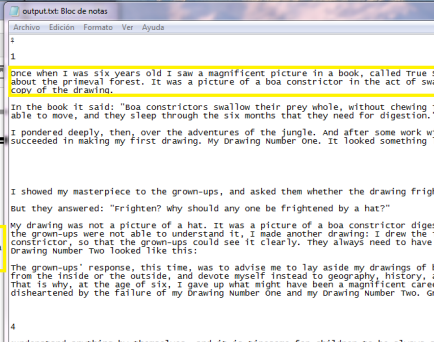
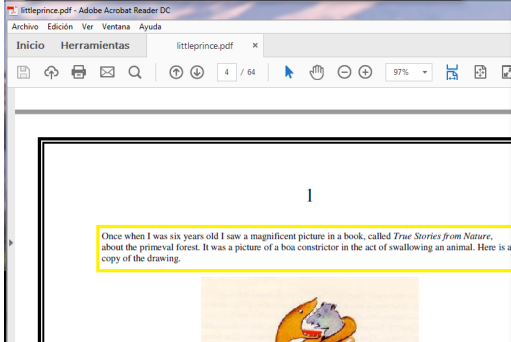
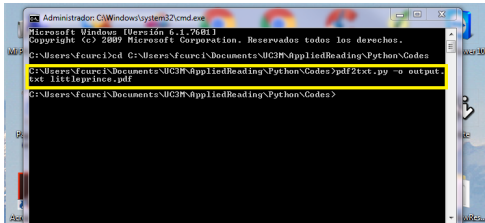
Scraping the little prince

The screenshot shows a Windows desktop environment. In the background, there is a wallpaper of a snow-capped mountain range. In the foreground, three windows are open:

- Administrador: C:\Windows\system32\cmd.exe**: A command prompt window showing the execution of a Python script. The command is `C:\Users\fcureci>cd C:\Users\fcureci\Documents\UC3M\AppliedReading\Python\Codes`, followed by `C:\Users\fcureci\Documents\UC3M\AppliedReading\Python\Codes>pdf2txt.py -o output.txt littleprince.pdf`, and finally `C:\Users\fcureci\Documents\UC3M\AppliedReading\Python\Codes>`.
- littleprince.pdf - Adobe Acrobat Reader DC**: A PDF viewer window displaying the first page of a document. The page has a large number '1' at the top. Below it, the text reads: "Once when I was six years old I saw a magnificent picture in a book, called *True Stories from Nature*, about the primeval forest. It was a picture of a boa constrictor in the act of swallowing an animal. Here is a copy of the drawing." Below the text is a small, colorful illustration of a yellow snake swallowing a grey mouse.
- output.txt: Bloc de notas**: A text editor window showing the output of the Python script. The text is as follows:
1
1
Once when I was six years old I saw a magnificent picture in a book, called *True Stories from Nature*, about the primeval forest. It was a picture of a boa constrictor in the act of swallowing an animal. Here is a copy of the drawing.
In the book it said: "Boa constrictors swallow their prey whole, without chewing." I was unable to move, and they sleep through the six months that they need for digestion.
I pondered deeply, then, over the adventures of the jungle. And after some work I succeeded in making my first drawing. My Drawing Number One. It looked something like this:

I showed my masterpiece to the grown-ups, and asked them whether the drawing frightened them. But they answered: "Frighten? why should any one be frightened by a hat?"
My drawing was not a picture of a hat. It was a picture of a boa constrictor digesting an elephant from the inside or the outside, and devote myself instead to geography, history, and science. That is why, at the age of six, I gave up what might have been a magnificent career as a painter. I am disheartened by the failure of my Drawing Number One and my Drawing Number Two. Grown-ups never understand anything by themselves, and it is tiresome for children to be always and forever understood by them.

Scraping the little prince



How Python can be useful to applied practitioners

- Fede C
 - What Python is and scraping
- Fede M
 - **GIS (geography and maps) application**
- Tomás
 - Machine learning tools

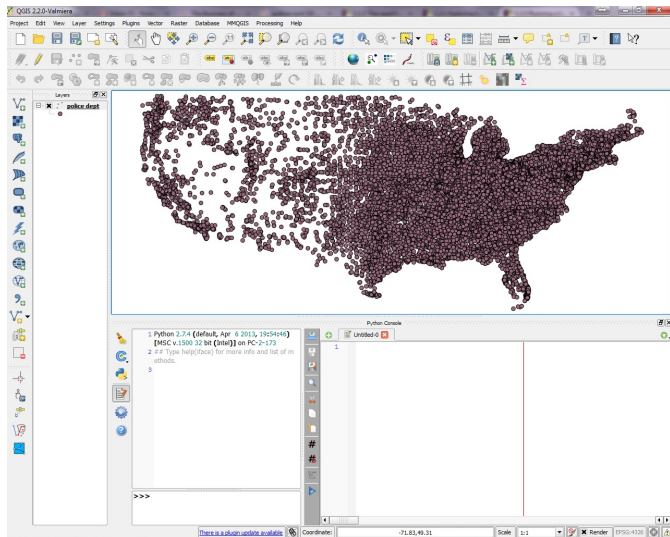
Python can be used both in QGIS and ArcGIS as a language to interact with these programs.

These will be useful to:

- Automated processes that you can already do by clicking in your preferred GIS program. (Very important for using data from STATA or creating data for STATA use)
- STATA can automate many of this operations already (not all of them) but GIS programs much more efficient
- Perform some new function only available if using Python

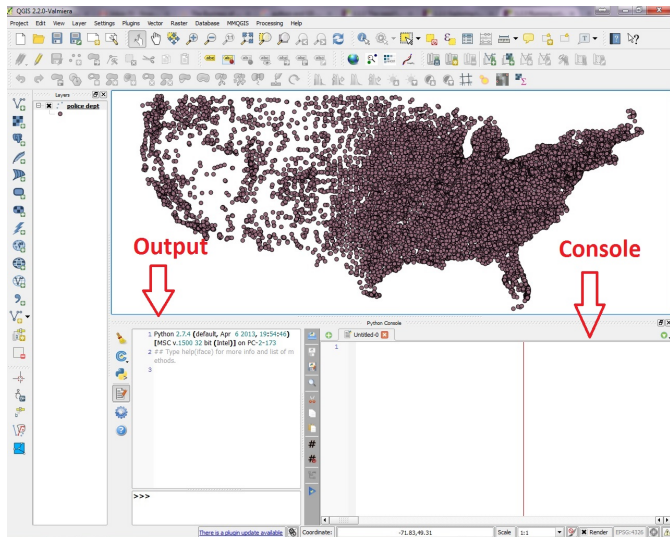
Where to Find Python? (QGIS)

Plugins ⇒ Python Console



Where to Find Python? (QGIS)

Plugins ⇒ Python Console



Learn to Write GIS Python Code

- 1 Follow the course found in *qgistutorials.com*
- 2 Search in *stackexchange.com*
- 3 Processing \Rightarrow “Graphical modeler” (a drag and drop tool that generates a Python Code)
- 4 Write your own .py file

Example 1 - Automation of Presentation of Results

You have shapefiles of all the crimes in counties in the US 1960-2014. For exploring the data you want to create a map for each year. Change the color of places with a high violent crime rate.

You can write your own .py code:

- 1 For loop for each file (for yy in year_available)
- 2 open .shp file - .addMapLayer()
- 3 select high crime places - .getFeatures()
- 4 Change color - .setSelectionColor()
- 5 save map

If you have multiple maps you can even created videos of maps

Example 2 - Automation Getting Inputs from STATA

Example: Data on number of police officers for all police departments and want to perform some operation with QGIS using this data. For QGIS to read data from STATA needs a .csv file with an identifier (ID of police department) and the variable of interest (number of police officers)

Python can easily automate all this in case you make some change to any variable you do not need to manually redo everything:

- 1 Open .shp file and .csv files - `.addMapLayer()`
- 2 Join the 2 files using the common identifier - `.addJoin`
- 3 Perform whatever operation you wish

Example 3 - Summary Statistics by Distance

Calculate number of police officers in a radius of 20km around a police department.

This could be done in STATA by calculating the distance of all police dept. from all police dept. and summing up the number of police officers if the distance is less than 20km. This is computationally unfeasible (15000 police dept.)

QGIS is very efficient when dealing with GIS data and distances. We can write a code on Python:

- 1 Open the file we created in example 2 - `.addMapLayer()`
- 2 Create a buffer of 20km around all police departments - `.buffer()`
- 3 Select only police departments that intersect with the buffer - `.intersects()`
- 4 Sum up police officers if it intersects and save in a new variable

Example 3 - Summary Statistics by Distance

Why is it useful?

- 1 Much faster than STATA (sometimes is just unfeasible on STATA)
- 2 Once you have the Python code is easy to change buffers or variables
- 3 You can then save the data and use this newly created variable back in STATA - `.writeAsVectorFormat()`

Example 4 - Neighbors of Polygons

All the previous examples the GIS data was points (police department). Many times data is in polygons (Counties, States,)

Let's say that we have population level data at the county level. We want to calculate the total population of the neighbors of all counties (where neighbors means to have one border in common). Python can do that.

Is a long code but the main tricky thing is how to determine neighbors. You do with the following function - `.intersects(geom.boundingBox())`

Again this new variable created in QGIS can be then exported in STATA

A Taste of More Complex Things you Can Do in Python

Excellent Guide: “Night Lights and ArcGIS: A Brief Guide” by Matt Lowe

Let's say you want to create a 20kmX20Km pixel-year panel of luminosity.

You can download a raster NASA NOAA satellite data of luminosity at night of the earth (data for many years and many types of measurements).

You can do it in Python with `.CreateFishnet_management()`

Given that the data is already on Python you can add a lot of information to this pixels dataset (Country, Distance to the coast,) and then save it in a readable file for STATA

How Python can be useful to applied practitioners

- Fede C
 - What Python is and scraping
- Fede M
 - GIS (geography and maps) application
- Tomás
 - **Machine learning tools**

Machine Learning with Scikit-Learn

- **Idea:** How to use a couple of methods of Imbens/Atthey NBER summer lectures with the library scikit-learn
- **Machine Learning...**
 - ...is more about fit/prediction.
 - ..cares a lot about scalability of the algorithm (big data, many covariates).
 - ...is neither about causality or inference (Susan Athey disagrees).
 - ...nor the formal asymptotic properties.
- **Focus on:** Out of sample performance, cross-validation, regularization.
- Why Python and Scikit-learn?

How to learn?

- Machine Learning:
 - **Coursera:** Stanford course with Andrew NG (matlab), Washington University specialization (python).
 - **edX:** Columbia course from the MicroMasters Program.
 - **Book:** *The Elements of Statistical Learning. Data mining, Inference and Prediction.* Trevor Hastie, Robert Tibshirani and Jerome Friedman.
- Basic packages to handle data on Python:
 - **Panda (Most popular library for data manipulation):** Resources.
 - **Numpy (Efficient python matrix handler):** Tutorial, lecture notes.
- Scikit-Learn:
 - **Links:** Official Tutorial, video lectures, other Tutorial .
 - **Book:** Learning scikit-learn: Machine Learning in Python.

Ex.1 Classification + Regularization

- **Built in scikit learn:** Lasso, Ridge regression, Least angle regression, Lars, Elastic Nets.
- **Example:** Reviews from an Amazon product. Which word has the highest predictive power of a positive review?
- **Method:** Regularized logistic regression.
- **Scikit implementation:** First, get your matrix of words using the *CountVectorizer*, then use *LogisticRegressionCV* with a Lasso like penalty ▶ `lasso`.
- *LogisticRegressionCV* has a built in cross validation to find the “optimal” penalty parameter α .

Ex.1 Classification + Regularization

- 166752 data points with 141224 covariates! “Only” around 3300 are nonzero.

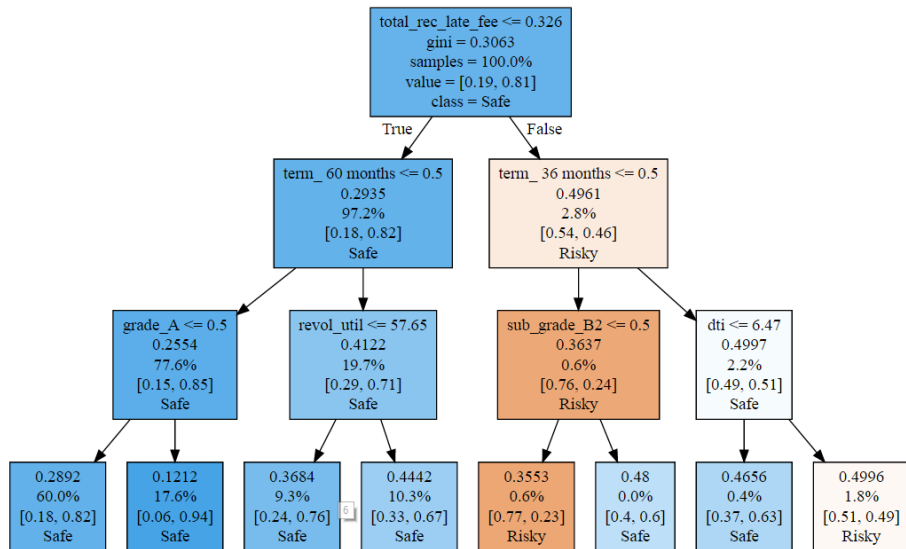
	Positive Words		Negative Words	
1	oustanding	2.12	dissapointed	-2.58
2	saves	2.03	worst	-2.47
3	lifesaver	1.90	worthless	-2.40
4	skeptical	1.80	theory	-2.38
5	adores	1.79	unusable	-2.33
6	south	1.69	disappointing	-2.26
7	con	1.67	rippoﬀ	-2.25
8	penny	1.65	poorly	-2.22
9	awesome	1.63	ineffective	-2.20
10	pleasantly	1.63	useless	-2.02

Ex.2 Regression Trees

- **Built in scikit learn:** Random Forests, Bagging, Boosting (Adaboost).
- **Example:** When to give a loan? Data on good/bad loans + covariates of the individual (employment, debt to income ratio...).
- **Method:** Random Forests: Fit different regression trees in random subsamples with replacement and apply a weighted average between them.
- **Scikit implementation:** *RandomForestClassifier*

```
forest =RandomForestClassifier(max_depth=4).fit(X,Y)
forest.predict_proba(X[1,:])
print forest.feature_importances_
```

Ex.2 Regression Trees



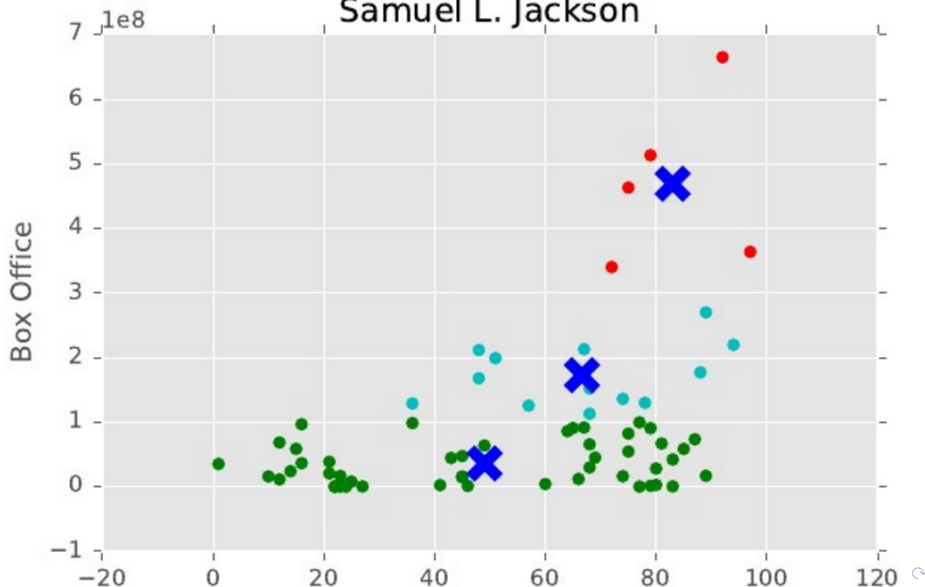
Ex.3 Clustering

- **Built in scikit learn:** K-means, Gaussian Mixtures.
- **Example:** Clustering Samuel L. Jackson movies by quality and revenue box office.
- **Method:** K-Means.
- **Scikit implementation:** *KMeans*

```
kmeans = KMeans(n_clusters=3,random_state=10).fit(data)
centroids = kmeans.cluster_centers_
labels = kmeans.labels_
```

Ex.3 Clustering

Samuel L. Jackson



Classifier + Regularization:

```
# Word Matrix
```

```
vectorizer = CountVectorizer(token_pattern=r'\b\w+\b')
```

```
words_matrix = vectorizer.fit_transform(products['review_clean'])
```

```
name = vectorizer.get_feature_names()
```

```
# Model
```

```
sentiment_model = linear_model.LogisticRegressionCV(penalty='l1',  
                                                    ..., solver='liblinear')
```

```
sentiment_model.fit(words_matrix, dependent_var)
```

```
coefficients = sentiment_model.coef_
```

▸ Lasso

$$\min_{\beta} \sum_i^N (Y_i - X_i \beta)^2 + \alpha ||\beta||_1 \quad (1)$$