# Tools for Management Data Science Course

Master in Big Data, Analytics and Technology Management
University of Firenze, 2016/2017
Andrea Gigli, PhD Applied Statistics, MSc Big Data & Social Mining

# Goals

Congrats, you enrolled in the **Master in Big Data, Analytics and Technology Management @ University of Firenze**, a Master organized by the Departments of Statistics & Computer Science, Engineering and Economics.

I'm Andrea Gigli, PhD in Applied Statistics, MSc in Big Data Analytics & Social Mining. During the Data Science for Management course I'll guide you through Data Science applications for Business.

In these **30 hours Front End lectures**, **90 hours Labs**, **60 hours Industrial Workshops** you'll understand selected business problems and learn how to use Data Science to solve them, while working on real cases and learning methodologies and techniques during your project works.

# Goals

Being in touch with real companies, I've designed this course in order to make your CV appealing for them. This year you'll learn

- R and Python scripting (R & Python)
- How to explore business data (R)
- How to build business dashboards using quantitative, qualitative and geo-referenced data (R)
- How to scrape web pages for business intelligence purposes (Python)
- How to measure customer engagement on Social Networks like facebook or twitter (Python)
- How to use data mining in your customer database (Python & R)
- How to segment and cluster customers (R)

Moreover, during the workshops you'll make practice with other problems and interact with true industry professionals to understand how to solve them with business in mind...

# Time to work now...

# Before coding…

I've collected some notes here on tools we'll use extensively during course, labs and seminars. You can see how to

1. install R and R Studio
2. install Python 3 and Spyder
3. install Jupyter in Linux Virtual Machine
4. get an access token for Facebook mining and how to use it
5. get an access token for Twitter mining
6. install Scrapy for Python 3.0
7. Xpath basics
8. GitHub basics

# Jupyter Notebook

To install IPython Notebook follow these steps

1.  Install jupyter for Python 3.X in your Linux

    ```
    $ sudo pip3 install jupyter
    ```
    (enter your password)

Now your Jupyter has Python (3) kernel. For using Jupyter with R you need to

2.  Install IRKernel in R

    ```
    install.packages(c('repr', 'IRdisplay', 'crayon', 'pbdZMQ', 'devtools'))

    devtools::install_github('IRkernel/IRkernel')

    IRkernel::installspec()  # to register the kernel in the current R installation
    ```
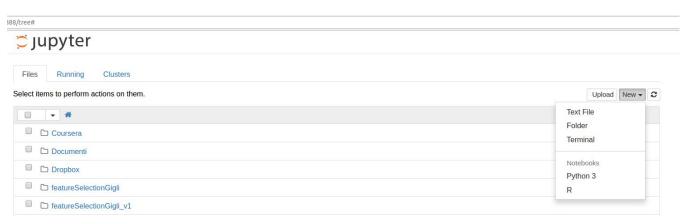
# Jupyter Notebook

To launch Jupyter from the Linux command window, type

3.  $ jupyter notebook

You'll see you browser opening a new session. If you click on "New" you can create a new Python or R notebook

# Facebook API

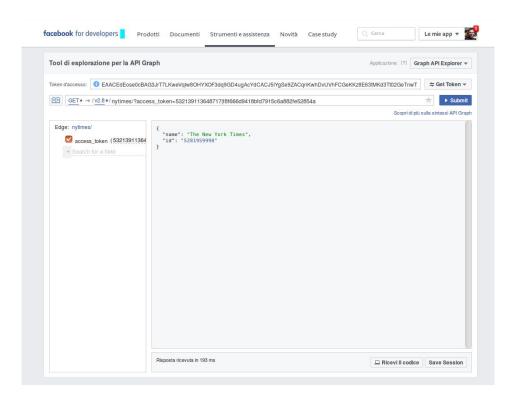To get access to Facebook API you need a Developer Account

1. Login to Facebook
2. Register to https://developers.facebook.com
3. Create a FB dev account and get a temporary (2 hours) token. You can also build a never expiring token by combining
   - app id: 123456789abcdef
   - secret key: f8f666d9418bfd7915c6a882fe52854a
   and separating them with a "pipe", "|"
   - never expiring token: 123456789abcdef|f8f666d9418bfd7915c6a882fe52854a

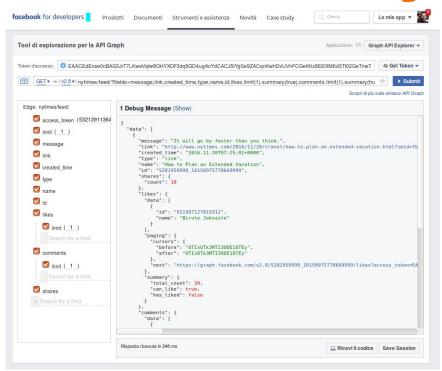You can find a very useful tutorial @ https://developers.facebook.com/apps

# Facebook API

To test your access token type

```
nytimes/?access_token =
<YOUR ACCESS TOKEN>
```
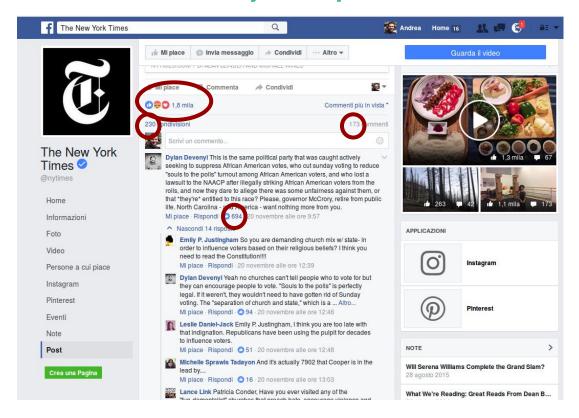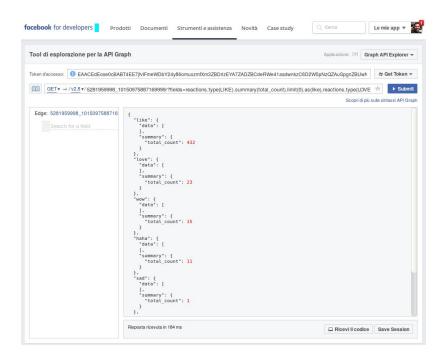
# Facebook API: messages



```
nytimes/feed/?fields=message
,link,created_time,type,name
,id,likes.limit(1).summary(t
rue),comments.limit(1).summa
ry(true),shares&limit=1&acce
ss_token=<YOUR ACCESS TOKEN>
```

# https://www.facebook.com/nytimes/posts/10150975887169999

# Facebook API: reactions

```
5281959998_10150975887169999/?field
s=reactions.type(LIKE).summary(tota
l_count).limit(0).as(like),reaction
s.type(LOVE).summary(total_count).l
imit(0).as(love),reactions.type(WOW
).summary(total_count).limit(0).as(
wow),reactions.type(HAHA).summary(t
otal_count).limit(0).as(haha),react
ions.type(SAD).summary(total_count)
.limit(0).as(sad),reactions.type(AN
GRY).summary(total_count).limit(0).
as(angry)&access_token=<YOUR ACCESS
TOKEN>
```
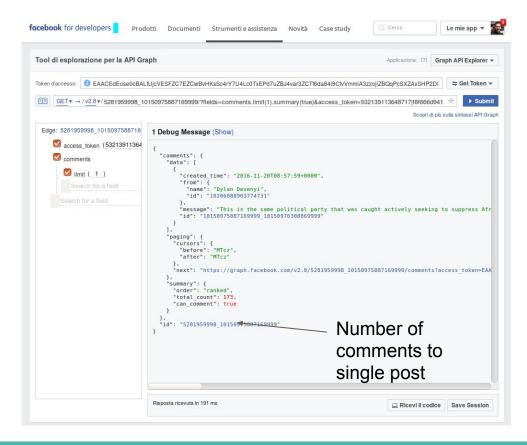
# Facebook API: comments to single post

```
5281959998_101509758871
69999/?fields=message,c
reated_time,comments.li
mit(1),summary(true)&ac
cess_token=532139113648
717|f8f666d9418bfd7915c
6a882fe52854a
```

See also

https://developers.facebook.co
m/docs/graph-api/reference/v2.
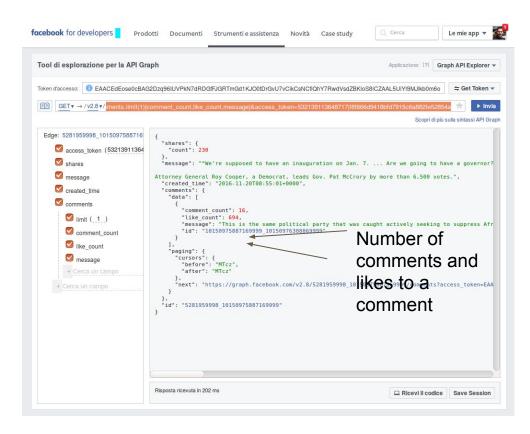8/comment/



Number of comments to single post

# Facebook API: comments counts to comments

5281959998_10150975887169
999/?fields=shares,messag
e,created_time,comments.l
imit(1){comment_count,lik
e_count,message}&access_t
oken=532139113648717|f8f6
66d9418bfd7915c6a882fe528
54a

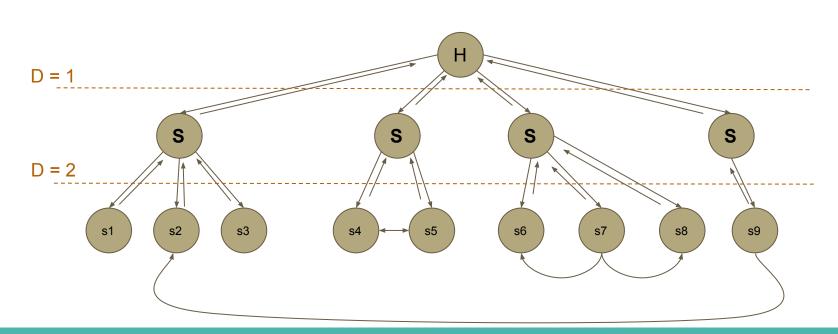See also

https://developers.facebook.com/d
ocs/graph-api/reference/v2.8/com
ment/

# Scrapy

Scrapy is a Python framework for building **web spider**, i.e. robots which are able to visit all the web pages they can reach through links.

# Scrapy

To install scrapy in Linux just type

```
sudo pip3 install scrapy
```

To be sure to have the most recent version type again

```
sudo pip3 install --upgrade scrapy
```

IMPORTANT! Sometimes Scrapy can generate conflicts when it is installed system-wide. We recommend to install it in a specific environment.

Look here >> https://doc.scrapy.org/en/latest/intro/install.html#intro-using-virtualenv

# Alternatives to Scrapy

Scrapy works well but sometimes you simply need **parsing wisely web pages** (I say wisely...) or to interact with javascript of **web pages**.

In these cases you should look at
- Beautiful Soup (web page parsing Python library, <u>start from here</u>)
- Selenium (tool for automating web application testing, <u>official Python doc here</u>)

Xpath logic is the same.

# XPath

Once you have downloaded the web pages you have to extract information from them.

**Xpath is a syntax** used to navigate through elements and attributes in an XML document.

XPath uses path expressions to select nodes or node-sets in an XML document.

These path expressions look very much like the expressions you use with a traditional computer file system ----------------------------------------------->

Folders
- BOOKS
  - BOOK
    - TITLE
    - AUTHOR
      - FIRSTNAME
      - LASTNAME

# XPath

In XPath, there are seven kinds of nodes: element, attribute, text, namespace, processing-instruction, comment, and document nodes.

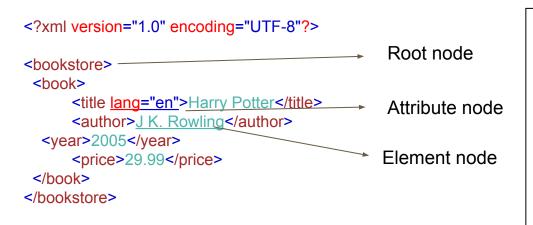XML documents are treated as trees of nodes. The topmost element of the tree is called the **root element** and **atomic values** are nodes with no children or parent.

```xml
<?xml version="1.0" encoding="UTF-8"?>

<bookstore>                                    Root node
 <book>
     <title lang="en">Harry Potter</title>     Attribute node
     <author>J K. Rowling</author>
  <year>2005</year>                            Element node
     <price>29.99</price>
 </book>
</bookstore>
```

For more details, please start with

http://www.w3schools.com/xml/xpath_intro.asp

and play with

http://codebeautify.org/Xpath-Tester

… and enjoy!

# GitHub

Git is **version control software**, which means it manages changes to a project without overwriting any part of that project.

This means you and your coworker(s) can each upload your revisions to the project (for example a web page, a paper, a program, …) and Git will save more copies. Later, you can merge your changes together without losing any work along the way. You can even revert to an earlier version at any time, because Git keeps a "snapshot" of every change ever made.

Git was designed with a big project like Linux in mind, there are a lot of Git commands. However, to use the basics of Git, you'll only need to know a few terms. They all begin the same way, with the word "**git**"

# GitHub: very frequent commands

| | |
|---|---|
| `git init` | Initializes a new Git repository. |
| `git clone` | Clone a Git repository and create a local repository |
| `git status` | Check the status of your repository. |
| `git pull` | "pull" the changes down from GitHub to your local repository. |
| `git add` | This does *not* add new files to your repository. Instead, it brings new files to Git's attention. |
| `git commit` | After you make any sort of change, you input this in order to take a "snapshot" of the repository. Usually it goes `git commit -m "Message here"` |
| `git push` | "push" the changes up to GitHub with this command. |
| `git help` | Forgot a command? Type this into the command line to bring up the 21 most common git commands. |

# GitHub

`git init:` Initializes a new Git repository. Until you run this command inside a repository or directory, it's just a regular folder. *Only after you input this does it accept further Git commands*.

`git config:` Short for "configure," this is most useful when you're setting up Git for the first time.

`git help:` Forgot a command? Type this into the command line to bring up the 21 most common git commands. You can also be more specific and type "git help init" or another term to figure out how to use and configure a specific git command.

`git status:` Check the status of your repository. See which files are inside it, which changes still need to be committed, and which branch of the repository you're currently working on.

`git add:` This does *not* add new files to your repository. Instead, it brings new files to Git's attention. After you add files, they're included in Git's "snapshots" of the repository.

`git branch:` Working with multiple collaborators and want to make changes on your own? *This command will let you build a new branch, or timeline of commits, of changes and file additions that are completely your own.* Your title goes after the command. If you wanted a new branch called "cats," you'd type `git branch cats`.

# GitHub

`git commit:` Git's most important command. After you make any sort of change, you input this in order to take a "snapshot" of the repository. Usually it goes `git commit -m "Message here."` The -m indicates that the following section of the command should be read as a message.

`git checkout:` Literally allows you to "check out" a repository that you are not currently inside. This is a navigational command that lets you move to the repository you want to check. You can use this command as `git checkout master` to look at the master branch, or `git checkout cats` to look at another branch.

`git merge:` When you're done working on a branch, you can *merge your changes back to the master branch*, which is visible to all collaborators. `git merge cats` would take all the changes you made to the "cats" branch and add them to the master.

`git push:` If you're working on your local computer, and *want your commits to be visible online on GitHub* as well, you "push" the changes up to GitHub with this command.

`git pull:` If you're working on your local computer and want *the most up-to-date version of your repository to work with*, you "pull" the changes down from GitHub with this command.

# GitHub

http://readwrite.com/2013/10/02/github-for-beginners-part-2/

# Regular Expression