Mineração de Dados Aula 1 – parte 3

Especialização em Ciência de Dados e suas Aplicações



Coleta







Permite obter dados facilmente

- Dois tipos básicos:
 - Streaming
 - Requisições (requisição-resposta)
- Disponíveis: Twitter, Foursquare, etc



Exemplo: Usando a API do Twitter

Passo1: Instalar a biblioteca TwitterAPI:

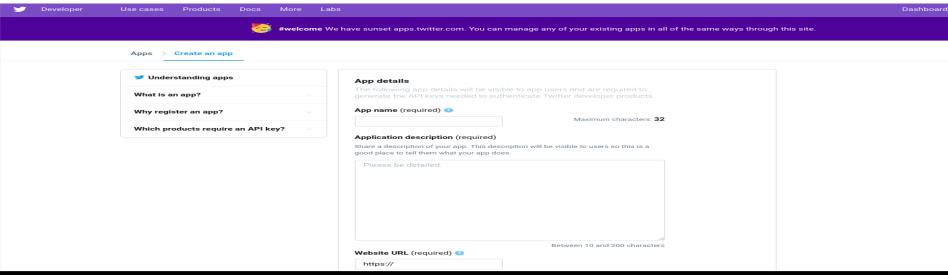
- pip install TwitterAPI
- ou fazer um clone no https://github.com/geduldig/TwitterAPI



Exemplo: Usando a API do Twitter

Passo 2. Criar uma aplicação no Twitter

- a. Necessário ter cadastro no Twitter
- **b.** Acessar https://developer.twitter.com/en/apps/create e clicar em "Create an App"





Exemplo: Usando a API do Twitter

Passo 3. Gerar os tokens da app

a. Na tela de detalhes da aplicação, clicar na aba "Keys and Access Token"

b. Para gerar o *token*, nesta mesma aba clique em "*Create Access Token*"

As opções podem ter mudado, mas a ideia continua a mesma.



Exemplo: Usando a API do Twitter

from TwitterAPI import TwitterAPI

Facilita a comunicação com a API do Twitter

```
twitter_api = TwitterAPI(consumer_key='XXX')
consumer_secret='XXXX',
access_token_key='XXXX',
access_token_secret='XXXX')
```

for item in twitter_api.request('statuses/filter', {'track':'csbc'}):
 print item['text']



Exemplo: Usando a API do Twitter

from TwitterAPI import TwitterAPI

Tokens obtidos

for item in twitter_api.request('statuses/filter', {'track':'csbc'}):
 print item['text']



Exemplo: Usando a API do Twitter

```
from TwitterAPI import TwitterAPI
```

for item in twitter_api.request('statuses/filter', {'track':'csbc'}):
 print item['text']

Requisição dos tweets com a termo "csbc"



Exemplo: Usando a API do Twitter

Demonstração - Coleta Streaming

coletaAPIs/streamingTwitter/exemplo1

coletaAPIs/streamingTwitter/exemplo2/coletaTweetsStreaming.py (Exemplo anterior com pequenas novidades)



Exemplo: Usando a API do Twitter

Exemplo de saída





Exemplo: Usando a API do Twitter



```
{"favorited": false, "contributors": null, "truncated": false, "text": "#CFP Workshop on Noisy User-generated Text at ACL - Beiji
ng 31 July 2015. Papers due: 11 May 2015. http://t.co/rcygyEowqH #NLProc_#WNUT15", "possibly sensitive": false, "in reply to st
atus id": null, "user": {"follow request sent": null, "profile use background image": true, "default profile image": false, "id":
237918251, "verified": false, "profile image url https": "https://pbs.twimg.com/profile images/527088456967544832/DnclpoZO norma
1.jpeg", "profile_sidebar_fill_color": "DDEEF6", "profile_text_color": "333333", "followers_count": 226, "profile_sidebar_borde
r color": "CODEED", "id str": "237918251", "profile background color": "CODEED", "listed count": 13, "profile background image ur
l_https": "https://abs.twimg.com/images/themes/theme1/bg.png", "utc_offset": null, "statuses_count": 120, "description": "I am a
postdoctoral researcher @PennCIS, studying Natural Larguage Processing and Social Media.", "friends count": 166, "location": "Phi
ladelphia PA", "profile link color": "0084B4", "profile image url": "http://pbs.twimg.com/profile images/527088456967544832/Dnclp
oZO normal.jpeg", "following": null, "geo enabled": true, "profile background image url": "http://abs.twimg.com/images/themes/the
me1/bg.png", "name": "Wei Xu", "lang": "en", "profile background tile": false, "favourites count": 88, "screen name": "cocoweix
u", "notifications": null, "url": "http://www.cis.upenn.edu/~xwe/", "created at": "Thu Jan 13 23:15:12 +0000 2011", "contributor
s enabled": false, "time zone": null, "protected": false, "default profile": true, "is translator": false}, "filter level": "lo
w", "geo": null, "id": 616333141884674048, "favorite_count": 0, "lang": "en", "entities": {"user_mentions": [], "symbols": [], "t
rends": [], "hashtags": [{"indices": [0, 4], "text": "CFP"}, {"indices": [124, 131], "text": "NLProc"}, {"indices": [132, 139],
"text": "WNUT15"}], "urls": [{"url": "http://t.co/rcygyEowqH", "indices": [99, 121], "expanded_url": "http://noisy-text.github.i
o", "display_url": "noisy-text.github.io"}]}, "in_reply_to_user_id_str": null, "retweeted": false, "coordinates": null, "timestam
p_ms": "1435780246598", "source": "<a href=\"http://twitter.com\" rel=\"nofollow\">Twitter Web Client</a>", "in_reply_to_status_i
d_str": null, "in_reply_to_screen_name": null, "id_str": "616333141884674048", "place": null, "retweet_count": 0, "created_at":
"Wed Jul 01 19:50:46 +0000 2015", "in reply to user id": null}
```



Formato JSON

Visualização da hierarquia



Outras opções para o Twitter, por exemplo:

http://www.tweepy.org/



Processando um JSON



Exemplo: JSON referente a um tweet

pip install simplejson

```
import simplejson as json
# Arquivo com Tweets
tweets file = open('tweet.txt', "r")
#le a linha do arquivo
tweet json = tweets file.readline()
#imprime a linha lida
print tweet json
                                        coletaAPIs/streamingTwitter/processaExemplo1/
#remove espacos em branco
strippedJson = tweet json.strip()
#converte uma string json em um objeto python
tweet = json.loads(strippedJson)
print tweet['id'] # ID do tweet
print tweet['created at'] # data de postagem
print tweet['text'] # texto do tweet
print tweet['user']['id'] # id do usuario que postou
print tweet['user']['name'] # nome do usuario
print tweet['user']['screen name'] # nome da conta do usuario
```



Exemplo: JSON referente a um tweet

```
import simplejson as json
# Arquivo com Tweets
tweets file = open('tweet.txt', "r")
#le a linha do arquivo
tweet json = tweets file.readline()
#imprime a linha lida
print tweet json
                                        coletaAPIs/streamingTwitter/processaExemplo1/
#remove espacos em branco
strippedJson = tweet json.strip()
#converte uma string ison em um objeto python
tweet = json.loads(strippedJson)
                                                              Acessando campos do
                                                              ison
print tweet['id'] # ID do tweet
print tweet['created at'] # data de postagem
print tweet['text'] # texto do tweet
print tweet['user']['id'] # id do usuario que postou
print tweet['user']['name'] # nome do usuario
print tweet['user']['screen name'] # nome da conta do usuario
```



Exemplo: JSON referente a um tweet

```
import simplejson as json
# Arquivo com Tweets
tweets file = open('tweet.txt', "r")
#le a linha do arquivo
tweet json = tweets file.readline()
#imprime a linha lida
print tweet json
#remove espacos em branco
strippedJson = tweet json.strip()
                                                          cocoweixu
#converte uma string json em um objeto python
tweet = json.loads(strippedJson)
print tweet['id'] # ID do tweet
print tweet['created at'] # data de postagem
print tweet['text'] # texto do tweet
print tweet['user']['id'] # id do usuario que postou
print tweet['user']['name'] # nome do usuario
print tweet['user']['screen name'] # nome da conta do usuario
```

```
616333141884674048
Wed Jul 01 19:50:46 +0000 2015
#CFP Workshop on Noisy User-generated Text
at ACL - Beijing 31 July 2015. Papers due:
11 May 2015. http://t.co/rcygyEowqH
#NLProc #WNUT15
237918251
Wei Xu
```



Exemplo: JSON referente a um tweet

```
import simplejson as json
# Arquivo com Tweets
tweets file = open('tweet.txt', "r")
#le a linha do arquivo
                                                        616333141884674048
tweet json = tweets file.readline()
                                                       Wed Jul 01 19:50:46 +0000 2015
                                                       #CFP Workshop on Noisy User-generated Text
#imprime a linha lida
                                                        at ACL - Beijing 31 July 2015. Papers due:
print tweet json
                                                        11 May 2015. http://t.co/rcygyEowqH
                                                       #NLProc #WNUT15
#remove espacos em branco
                                                        237918251
strippedJson = tweet json.strip()
                                                        Wei Xu
                                                        cocoweixu
#converte uma string json em um objeto python
tweet = json.loads(strippedJson)
                                     Processando todos os tweets retornados em
print tweet['id'] # ID do tweet
print tweet['created at'] # data de postagem
                                                  saidaColetaStream.txt.
print tweet['text'] # texto do tweet
print tweet['user']['id'] # id do usualictatePis/streamingTwitter/processaExemplo2/
print tweet['user']['screen name'] # nome da conta do usuario
```



Personalização da coleta

```
r = api.request('statuses/filter', {'locations':'-74,40,-73,41'})
for item in r:
    print item
```

Longitude e latitude da área: NYC Tweets geolocalizados nessa área

Parameter value	Tracks Tweets from
-122.75,36.8,-121.75,37.8	San Francisco
-74,40,-73,41	New York City
-122.75,36.8,-121.75,37.8,-74,40,-73,41	San Francisco OR New York City



Personalização da coleta

```
r = api.request('statuses/filter', {'locations':'-74,40,-73,41'})
for item in r:
    print item
```

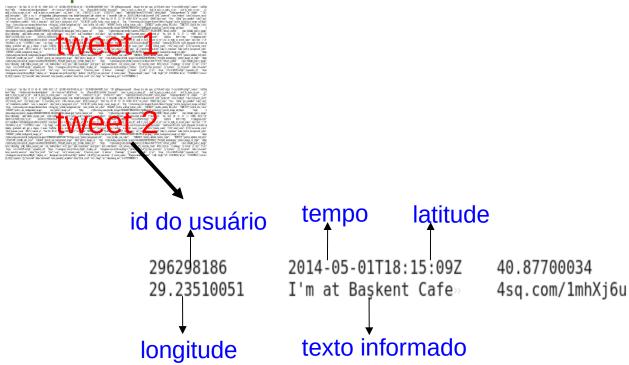
Longitude e latitude da área: NYC Tweets geolocalizados nessa área

Parameter value	Tracks Tweets fro	n Combinação
-122.75,36.8,-121.75,37.8	San Francisco	track=twitter&locations=-122.75,36.8,-121.75,37.8
-74,40,-73,41	New York City	Dolovro "twitter" au twoote de
-122.75,36.8,-121.75,37.8,-74,40,-73,41	San Francisco OR New York City	Palavra "twitter" ou tweets de San Fran

Assim será retornado tweets fora da área com a palavra Twitter





















Web crawler

Programas que analisam páginas Web em busca do conteúdo desejado

- · Alternativa para quando não existe API
- · Existem várias estratégias



```
# -*- encoding: utf-8 -*-
import urllib

#url obtida através do tweet
url = "http://www.utfpr.edu.br"

pagina = urllib.urlopen(url).read()

print(pagina)

fsaida = open('paginaColetada.html','w')
fsaida.write(str(pagina))
```

Demonstração coletaCrawler/ColetaPaginaSimples.py

-*- encoding: utf-8 -*-



```
import urllib
  #url obtida através do tweet
  url = "http://www.utfpr.edu.br"
  pagina = urllib.urlopen(url).read()
fsaida = open('paginaColetauu...
fsaida.write(str(pagina))

rt" class="primeix on the pagina of the 
  print(pagina)
                                                                                                                           cp"><a href="http://www.utfpr.edu.br/cornelioprocopio">Cornélio Procópio</a></
                                                                                                                           fb"><a href="http://www.utfpr.edu.br/franciscobeltrao">Francisco Beltrão</a></
                                                                                                                           qp" class="primeiro"><a href="http://www.utfpr.edu.br/quarapuava">Guarapuava</
                                                                                                                           ld"><a href="http://www.utfpr.edu.br/londrina">Londrina</a>
                                                                                                                           md"><a href="http://www.utfpr.edu.br/medianeira">Medianeira</a>
                                                                                                                           pb"><a href="http://www.utfpr.edu.br/patobranco">Pato Branco</a>
                                                                                                                           pg"><a href="http://www.utfpr.edu.br/pontagrossa">Ponta Grossa</a>
                                                                                                                           sh"><a href="http://www.utfpr.edu.br/santahelena">Santa Helena</a>
                                                                                                                           td"><a href="http://www.utfpr.edu.br/toledo">Toledo</a>
```



Exemplo com a biblioteca BeautifulSoup

BeautifulSoup para copiar dados de páginas Web

Para instalá-la, execute o comando:

pip install beautifulsoup4 sudo apt-get install python-bs4



```
[Bom dia, Fafich!]
Bom dia, Fafich!
```

coletaCrawler/beautifulSoup/example2.py

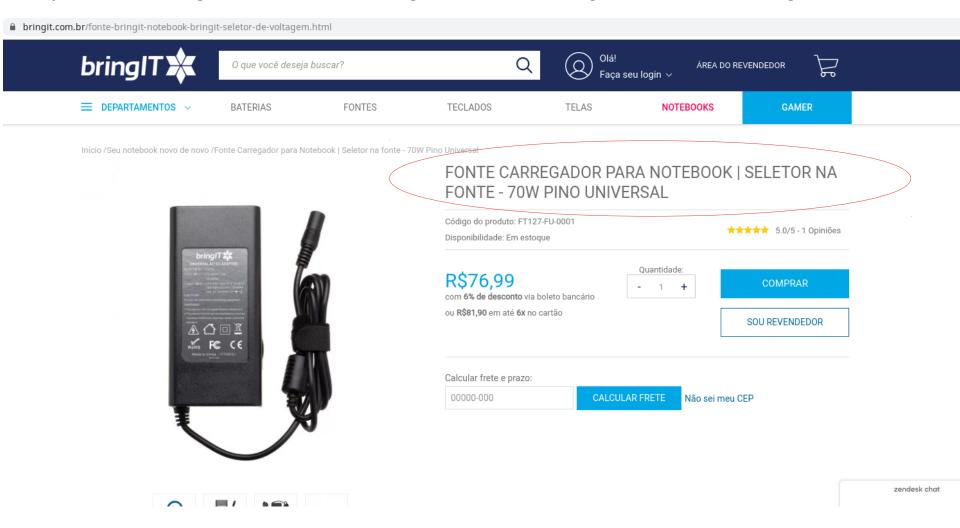


```
from bs4 import BeautifulSoup
pagina = "onetwo"
soup = BeautifulSoup(pagina)
print allTR _____
allTD = soup.find_all('td')
print allTD_____
                       → [one, two]
for t in soup.find_all('td'):
 print t ______ one</rr>
                        two
for t in soup.find_all('td'):
 print t.text ____
                        two
```

coletaCrawler/beautifulSoup/example3.py



https://www.bringit.com.br/fonte-bringit-notebook-bringit-seletor-de-voltagem.html



Pegar o nome do produto



```
import requests
from bs4 import BeautifulSoup
##pagina do produto
url = "https://www.bringit.com.br/fonte-bringit-notebook-bringit-seletor-de-voltagem.html"
##retorna o conteudo da pagina
req = requests.get(url)
##transforma o conteudo da pagina em um objeto BeautifulSoup
soup = BeautifulSoup(reg.content, 'html.parser')
nomeBruto = soup.find("div",{"class":"product-name"})
print nomeBruto.text
```

Faça um código para pegar o preço

coletaCrawler/beautifulSoup/bringit/



Algumas informações podem não ser coletáveis com o Beautifulsoup

O uso do Selenium pode ser demandado

Ver exemplo em: coletaCrawler/selenium/exemploSelenium.py



```
from selenium import webdriver
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as ec
from selenium.webdriver.common.by import By
from bs4 import BeautifulSoup
driver = webdriver.Chrome()
driver.get("https://www.w3schools.com/xml/ajax_intro.asp")
#Clica no botao da pagina para ver o conteudo
botaolist = driver.find_elements_by_xpath('//*[@id="demo"]/button')
botaolist[0].click()
wait = WebDriverWait(driver, 10)
wait.until(ec.visibility_of_element_located((By.XPATH, '//*[@id="demo"]/h1')))
# get the page source
page_source = driver.page_source
#fecha o driver, mas quando estivermos coletando varias paginas podemos manter ativo pra nao precisar abrir o browser novamente
#driver.close()
soup = BeautifulSoup(page_source, "lxml") #grab the content with beautifulsoup for parsing
content= soup.find("div",{"id":"demo"})
print content
```

Thiago H. Silva