Raspagem de Dados para mulheres



26 de novembro, das 9h às 17h

Inscrições pelo sympla.com.br/PyLadiesSP

Apoio:















Raspagem de dados ou web scraping é a coleta automatizada de dados na Internet.

Há uma variedade de técnicas e tecnologias para fazer isso.

- Para coleta de dados não estruturados como HTML que misturam dados e apresentação:
 - Uso de bibliotecas para extrair os dados do HTML (BeautifulSoup, Scrapy, ...)
- Para coleta de dados estruturados em formato XML ou JSON:
 - Uso de APIs disponibilizadas pelos provedores dos dados
 - Uso de biblioteca para manipular json ou xml





- 1. armazenar em arquivos, banco de dados, ...
- 2. analizar
- 3. apresentar (graficamente)

HTML



(escreva em um arquivo: ws_html.py)

from urllib.request import urlopen



```
# url de algum site
url = 'http://www.bbc.com/mundo'
# acessando a url e atribuindo a uma variável
# response é do tipo http.client.HTTPResponse (<http.client.HTTPResponse object at
0x0316F470>)
response = urlopen(url)
# resultado atribuido a resultado
resultado = response.read().decode('utf-8')
# gravar o resultado em um arquivo
open('resultado_html.txt', 'w', encoding='utf-8').write(resultado)
# gravar o resultado em um arquivo
open('resultado_html.html', 'w', encoding='utf-8').write(resultado)
```

Analisar o Resultado da Coleta de HTML



- Abra o arquivo resultado_html.html
- Abra o arquivo resultado_html.txt

O resultado_html.html mostrará o conteúdo visual do que foi coletado, pois o HTML é um formato de apresentação.

O resultado_html.txt mostrará o código HTML do mesmo conteúdo.

O resultado não contém dados estruturados, pois o HTML mistura dados + apresentação.

Desta forma fica difícil extrair os dados.





biblioteca para manipular o código HTML

é necessário instalá-la:

pip install beautifulsoup4

Como importar:

from bs4 import BeautifulSoup





```
from urllib.request import urlopen
from bs4 import BeautifulSoup

html = urlopen('http://www.python.org/')
bs0bj =
BeautifulSoup(html.read(),'html.parser')
print(bs0bj.h1)
```

```
$ terminal: python arquivo.py
<h1 class="site-headline">
<a href="/"><img alt="python<sup>TM</sup>" class="python-logo"
src="/static/img/python-logo.png"/></a>
</h1>
```





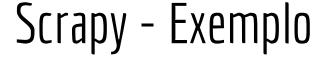
"[...] biblioteca Python que manipula grande parte da complexidade de busca e avaliação dos links de um site, rastreando domínios ou lista de domínios com facilidade."

Para instalar:

pip install scrapy

Como importar:

import scrapy





```
import scrapy
class SpiderSimples(scrapy.Spider):
    name = 'meuspider'
    start urls = 'http://example.com'
    def parse(self, response):
         self.log('Visitei o site: %s' % response.url)
         yield 'url': response.url, 'tamanho' len(response.body)}
         proxima url = 'http://www.google.com.br'
         self.log('Agora vou para: %s' % proxima url)
         yield scrapy.Request(proxima_url, self.handle_google)
    def handle_google(self, response):
         self.log('Visitei o google via URL: %s' % response.url)
```

Scrapy - Exemplo



Para rodar, utilize o comando:

```
$ terminal:
scrapy runspider arquivo.py
```

RSS (XML)



(escreva em um arquivo: ws_rss.py)

from urllib.request import urlopen



```
# url de algum site
url = 'http://feeds.bbci.co.uk/news/world/rss.xml'
# acessando a url e atribuindo a uma variável
# response é do tipo http.client.HTTPResponse (<http.client.HTTPResponse object at
0x0316F470>)
response = urlopen(url)
# resultado atribuido a resultado
resultado = response.read().decode('utf-8')
# gravar o resultado em um arquivo
open('resultado_rss.txt', 'w', encoding='utf-8').write(resultado)
# gravar o resultado em um arquivo
open('resultado rss.xml', 'w', encoding='utf-8').write(resultado)
```





- Abra o arquivo resultado_rss.xml
- Abra o arquivo resultado_rss.txt

O resultado_rss.xml mostrará o conteúdo visual do que foi coletado, pois o rss.xml é um padrão que pode ser "renderizado" pelo navegador.

O resultado_rss.txt mostrará o código rss.xml do mesmo conteúdo.

O resultado contém dados estruturados no padrão rss.



<?xml version="1.0" encoding="UTF-8"?> <?xml-stylesheet title="XSL formatting" type="text/xsl" href="/shared/bsp/xsl/rss/nolsol.xsl"?> <rss xmlns:dc="http://purl.org/dc/elements/1.1/" xmlns:content="http://purl.org/rss/1.0/modules/content/" xmlns:atom="http://www.w3.org/2005/Atom" version="2.0" x0</pre> <channel> <title><! [CDATA[BBC News - World]]></title> <description><! [CDATA[BBC News - World]]></description> <link>http://www.bbc.co.uk/news/</link> <image> <url>http://news.bbcimg.co.uk/nol/shared/img/bbc news 120x60.gif</url> <title>BBC News - World</title> <link>http://www.bbc.co.uk/news/</link> </image> <generator>RSS for Node</generator> <lastBuildDate>Sat, 19 Nov 2016 19:43:07 GMT</lastBuildDate> <copyright><![CDATA[Copyright: (C) British Broadcasting Corporation, see http://news.bbc.co.uk/2/hi/help/rss/4498287.stm for terms and conditions of reuse <language><! [CDATA[en-gb]]></language> <ttl>15</ttl> <item> <title><![CDATA[Syria conflict: Aleppo hospitals 'knocked out by bombardment']]></title> <description><! [CDATA [Days of strikes on Aleppo leave rebel areas virtually without a functioning hospital, reports say.]]></description> <link>http://www.bbc.co.uk/news/world-middle-east-38039282</link> <quid isPermaLink="true">http://www.bbc.co.uk/news/world-middle-east-38039282 <pubDate>Sat, 19 Nov 2016 19:10:20 GMT</pubDate> <media:thumbnail width="976" height="549" url="http://c.files.bbci.co.uk/3166/production/ 92564621 036476169-1.jpg"/> </item> <item> <title><![CDATA[Pope Francis names 17 new cardinals of Roman Catholic Church]]></title> <description><![CDATA[Pope Francis names 17 new cardinals from around the world, many of whom will help name his successor.]]></description> http://www.bbc.co.uk/news/world-europe-38038501</link> <quid isPermaLink="true">http://www.bbc.co.uk/news/world-europe-38038501 <pubDate>Sat, 19 Nov 2016 13:00:42 GMT</pubDate> <media:thumbnail width="976" height="549" url="http://c.files.bbci.co.uk/1350C/production/ 92561197 mediaitem92561195.jpg"/> </item> <item> <title><!!CDATA|Trump locks horns with Hamilton musical cast over Pence booing||></title> <description><!(CDATA(Donald Trump wades into a row over the booing of the US vice-president elect at a Broadway musical.))></description> <link>http://www.bbc.co.uk/news/world-us-canada-38039286</link> <quid isPermaLink="true">http://www.bbc.co.uk/news/world-us-canada-38039286</quid>





Temos dois arquivos com conteúdo muito parecido.

Podemos melhorar o código.

- transformar as instruções em funções
- juntar em um módulo o que é em comum entre os dois arquivos
- separar o que é diferente

```
# arquivo com instruções em comum: ws.py
```



from urllib.request import urlopen

```
def coletar dados(url):
   # acessando a url e atribuindo a um variável
   # response é do tipo http.client.HTTPResponse
(<http.client.HTTPResponse object at 0x0316F470>)
   response = urlopen(url)
   # resultado atribuido a resultado
   return response.read().decode('utf-8')
def guardar dados(nome arquivo, resultado):
   # gravar o resultado em um arquivo
   open(nome arquivo, 'w', encoding='utf-8').write(resultado)
```



```
# arquivo com instruções em comum: ws.py (parte 2)
def atualizar_arquivo(nome_arquivo, resultado):
    # gravar o resultado em um arquivo
    open(nome_arquivo, 'a', encoding='utf-8').write(resultado)
```



```
# arquivo com instruções específicas: ws_html.py
import ws

url = 'http://www.bbc.com/mundo'

dados = ws.coletar_dados(url)
ws.guardar_dados('resultado_html.html', dados)
```



```
# arquivo com instruções específicas: ws rss.py
import ws
url =
'http://feeds.bbci.co.uk/news/world/rss.xml'
dados = ws.coletar dados(url)
ws.guardar dados('resultado rss.html', dados)
```





Os fornecedores dos dados disponibilizam APIs (Application Programming Interfaces) para facilitar a coleta dos dados.

Geralmente elas requerem autenticação, ou seja, uma API key e access token.

Dados abertos: educação no Brasil



Exemplo: http://educacao.dadosabertosbr.com/api

http://educacao.dadosabertosbr.com/api/docs/%2Fapi%2Fescolas%2Fbuscaavancada

http://educacao.dadosabertosbr.com/api/escolas/buscaavancada?situacaoFuncionamento=1

```
[ ▼ 2 items, 34 KB
  195656.
  [ ▼ 100 items, 34 KB
    { ▼ 14 properties, 357 bytes
      "anoCenso": 2013.
      "cod": 31315001.
      "nome": "-ASSOCOACAO ASSISTENCIAL CRECHE JARDIM SANTA CECILIA",
      "codCidade": 3136702.
      "cidade": "JUIZ DE FORA",
      "estado": "MG".
      "regiao": "Sudeste".
      "situacaoFuncionamento": 1,
      "dependenciaAdministrativa": 4,
      "idebAI": 0.0.
      "idebAF": 0.0.
      "enemMediaGeral": 0.0.
      "dependenciaAdministrativaTxt": "Privada",
      "situacaoFuncionamentoTxt": "Em atividade"
    },
    { ▼ 14 properties, 371 bytes
      "anoCenso": 2013.
      "cod": 33062501.
      "nome": "0101001 ESCOLA MUNICIPAL VICENTE LICINIO CARDOSO",
      "codCidade": 3304557.
      "cidade": "RIO DE JANEIRO".
      "estado": "RJ",
      "regiao": "Sudeste".
      "situacaoFuncionamento": 1,
      "dependenciaAdministrativa": 3,
      "idebAI": 0.0,
      "idebAF": 4.900000095367432.
      "enemMediaGeral": 0.0,
      "dependenciaAdministrativaTxt": "Municipal",
      "situacaoFuncionamentoTxt": "Em atividade"
```





```
import csv

url =
'http://educacao.dadosabertosbr.com/api/escolas/buscaavancada?s
ituacaoFuncionamento=1'
dados = ws.coletar_dados(url)
dados_escolas = json.loads(dados)
quantidade = dados_escolas [0]
escolas = dados_escolas [1]
```



```
import csv
with open('escolas.csv', 'w', newline= '') as csvfile:
    fieldnames = ['nome', 'cidade', 'estado',
'dependenciaAdministrativa', 'enemMediaGeral']
    writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
    writer.writeheader()
    for escola in escolas:
        spamwriter.writerow(escola)
```

Facebook API - Preparação



https://developers.facebook.com/

Entrar

Adicionar um aplicativo

Configurar o aplicativo

Obter um ID do aplicativo

Sobre a API:

https://developers.facebook.com/docs/graph-api





















Aplicativo nativo ou para desktop?

Ative se seu aplicativo for nativo ou um aplicativo para desktop





Script para:

- obtenção do token de acesso
- consulta de posts

```
# arquivo instruções específicas: ws_facebook_v1.py (parte 1)
import ws
# configurado no facebook
app id = '1522873017729918'
app secret = 'b051937827b5a8de756121263bd76f08'
url para obter access token =
'https://graph.facebook.com/oauth/access token?client id={app id
}&client_secret={app_secret}&grant_type=client_credentials'.form
at(
    app id=app id,
    app secret=app secret)
# coletar os dados
access_token = ws.coletar_dados(url_para_obter_access token)
print(access token)
```



```
# arquivo instruções específicas: ws_facebook_v1.py (parte 2)
api_url = 'https://graph.facebook.com/'
page_id = 'PyLadiesSP'
url = api_url + page_id + '?fields=posts&' + access_token

print(url)
dados = ws.coletar_dados(url)
ws.guardar_dados('resultado_facebook.json', dados)
```





- criar uma função para a obtenção do token
- criar uma função para gerar a url para consultar posts de qualquer usuário do Facebook

```
# arquivo ws_facebook_v2.py (parte 1)
import ws
def obter access token():
    # configurado no facebook
    app id = '1522873017729918'
    app secret = 'b051937827b5a8de756121263bd76f08'
    url para obter access token =
'https://graph.facebook.com/oauth/access token?client id={app id
}&client secret={app secret}&grant type=client credentials'.form
at(
```

return ws.coletar_dados(url_para_obter_access_token)

app id=app id,

app secret=app secret)

coletar os dados

```
# ws facebook v2.py (parte 2)
api_url = 'https://graph.facebook.com/'
def url para obter posts(page id, access token):
    return api_url + page_id + '?fields=posts&' + access token
access token = obter access token()
# posts de PyLadiesSP
url = url para obter posts('PyLadiesSP', access token)
dados = ws.coletar dados(url)
ws.guardar dados('resultado facebook pyladiessp.json', dados)
# posts de Ideias Incríveis
url = url_para_obter_posts('ideias.incriveis', access token)
dados = ws.coletar_dados(url)
ws.guardar dados('resultado facebook ideiasincriveis.json',
dados)
```

JSON



```
{ ▼ 2 properties, 11 KB
 "posts": { V 2 properties, 11 KB
   "data": [ V 25 items, 11 KB
     { ▼ 4 properties, 297 bytes
        "message": "As meninas do Pyladies dominaram a Vila das Minas na Roadsec esse ano com workshops, dojo e palestras \u003C3\n#souPyLadiesSP #roadsec",
        "story": "PyLadies S\u00e3o Paulo shared PyLadies Brasil's album.",
        "created time": "2016-11-18T22:22:35+0000",
        "id": "1469656506672137 1615174665453653"
      { ▼ 4 properties, 310 bytes
        "message": "J\u00e1 come\u00e7ou o Workshop de Introdu\u00e7\u00e3o \u00e0 Python aqui na Roadsec.\nE s\u00f3 chegar aqui na \u00e1rea de atividades.
        ;) \n\n#souPvLadiesSP #roadsec",
        "story": "PyLadies S\u00e3o Paulo at Roadsec.".
        "created time": "2016-11-18T13:09:35+0000",
        "id": "1469656506672137 1615021138802339"
      { ▼ 4 properties, 661 bytes
        "message": "Amanh\u00e3 tem Roadsec e estamos com v\u00e1rias atividades na Vila das Minas. Confira a programa\u00e7\u00e3o:\n10h - Workshop de
        introdu\u00e7\u00e3o a Python com as PyLadies S\u00e3o Paulo\n11h - Mesa Redonda: Quem s\u00e3o as mulheres na computa\u00e7\u00e3o com as PyLadies
        Brasil\n13h - Dojo com as PvLadies S\u00e3o Carlos\n15h - Workshop de introdu\u00e7\u00e3o a Machine Learning com as PvLadies Campinas\n15h - Palestra
        PyLadies S\u00e3o Paulo e o que \u00e9 Python\nE ai quem vamos? \u003C3",
        "story": "PyLadies S\u00e3o Paulo shared Roadsec's photo.".
        "created time": "2016-11-18T01:21:43+0000",
        "id": "1469656506672137 1614868388817614"
      { ▼ 3 properties, 293 bytes
        "message": "O Grupy-SP est\u00e1 completando 9 anos com um mega meetup!!!\nSer\u00e1 dia 19\/11 e tem muito espa\u00e7o pra todo mundo poder ir.
        Inscri\u00e7\u00f5es no link abaixo.\nhttp:\/\/meetu.ps\/31YNpm",
        "created time": "2016-11-10T23:54:47+0000".
        "id": "1469656506672137 1612428422394944"
      { ▼ 4 properties, 280 bytes
        "message": "Estaremos na Vila das Minas da Roadsec dia 18//11 com um Workshop de introdu\u00e7\u00e3o a Python.\nE ai guem vamos? \u003C3",
        "story": "PvLadies S\u00e3o Paulo shared Roadsec's photo.".
        "created time": "2016-11-09T17:36:01+0000",
        "id": "1469656506672137 1611899782447808"
```

JSON



json.loads converte uma string que contém formato JSON em dicionário.

Assim podemos manipular os dados.

```
import json
json.loads('["foo", {"bar":["baz", null, 1.0, 2]}]')
['foo', {'bar': ['baz', None, 1.0, 2]}]
```



```
print(dados['posts']['data'][0]['message'])
As meninas do Pyladies dominaram a Vila da ...
```

```
print(dados['posts']['data'][1]['id'])
1469656506672137_1615021138802339
```





- incluir no início import json
- converter json para dict
- manipular os dados no dicionário
- criar um arquivo que contenha os textos dos posts

```
# arquivo ws facebook v3.py (parte 1)
import son
import ws
def obter access token():
    # configurado no facebook
    app id = '1522873017729918'
    app secret = 'b051937827b5a8de756121263bd76f08'
    url para obter access token =
'https://graph.facebook.com/oauth/access token?client id={app id
}&client_secret={app_secret}&grant_type=client_credentials'.form
at(
    app id=app id,
    app secret=app secret)
    # coletar os dados
    return ws.coletar_dados(url_para_obter_access_token)
```

```
# ws facebook v3.py (parte 2)
api_url = 'https://graph.facebook.com/'
def url para obter posts(page id, access token):
    return api url + page id + '?fields=posts&' + access token
access token = obter access token()
# posts de PyLadiesSP
url = url para obter posts('PyLadiesSP', access token)
dados = ws.coletar dados(url)
ws.guardar dados('resultado facebook pyladiessp.json', dados)
# converte json em dict
posts = json.loads(dados)
```



```
# ws_facebook_v3.py (parte 3)

messages = []

for post in posts['posts']['data']:
    messages.append(post.get('message', post.get('story', ''))))

ws.guardar_dados('posts.txt', '\n'.join(messages))
```

```
{ ▼ 2 properties, 11 KB
 "posts": { ▼ 2 properties, 11 KB
   "data": [ ▼ 25 items, 11 KB
     { ▼ 4 properties, 297 bytes
        "message": "As meninas do Pyladies dominaram a Vila das Minas na Roadsec esse ano com workshops, dojo e palestras \u003C3\n#souPyLadiesSP #roadsec",
        "story": "PyLadies S\u00e3o Paulo shared PyLadies Brasil's album.",
        "created time": "2016-11-18T22:22:35+0000",
        "id": "1469656506672137 1615174665453653"
      { ▼ 4 properties, 310 bytes
        "message": "J\u00e1 come\u00e7ou o Workshop de Introdu\u00e7\u00e3o \u00e0 Python agui na Roadsec.\nE s\u00f3 chegar agui na \u00e1rea de atividades.
        ;) \n\n#souPyLadiesSP #roadsec",
        "story": "PyLadies S\u00e3o Paulo at Roadsec.",
        "created time": "2016-11-18T13:09:35+0000",
        "id": "1469656506672137 1615021138802339"
      { ▼ 4 properties, 661 bytes
        "message": "Amanh\u00e3 tem Roadsec e estamos com v\u00e1rias atividades na Vila das Minas. Confira a programa\u00e7\u00e3o:\n10h - Workshop de
        introdu\u00e7\u00e3o a Python com as PyLadies S\u00e3o Paulo\n11h - Mesa Redonda: Quem s\u00e3o as mulheres na computa\u00e7\u00e3o com as PyLadies
        Brasil\n13h - Dojo com as PyLadies S\u00e3o Carlos\n15h - Workshop de introdu\u00e7\u00e3o a Machine Learning com as PyLadies Campinas\n15h - Palestra
```





• criar uma função para obter apenas as mensagens dos posts



```
# ws facebook_v4.py (parte 2)
api url = 'https://graph.facebook.com/'
def url para obter posts(page id, access token):
    return api_url + page_id + '?fields=posts&' + access token
def obter mensagens(posts):
    messages =
    for post in posts['posts']['data']:
       messages.append(post.get('message', post.get('story',
'')))
    return messages
```

```
# ws facebook_v4.py (parte 3)
access token = obter access token()
# posts de PyLadiesSP
url = url para obter posts('PyLadiesSP', access token)
dados = ws.coletar dados(url)
posts = json.loads(dados)
mensagens = obter mensagens(posts)
ws.guardar dados('posts.txt', '\n'.join(mensagens))
```







A cada consulta de uma API, além dos dados, é também informado quais são os endereços para buscar os resultados anteriores e os posteriores, como se fosse a paginação dos dados.

ws_facebook_v5.py



repetidas consultas



```
{ ▼ 2 properties, 11 KB
  "posts": { V 2 properties, 11 KB
    "data": [ ▼ 25 items, 11 KB

  ▼ 4 properties, 297 bytes

         "message": "As meninas do Pyladies dominaram a Vila das Minas na Roadsec esse ano com workshops, dojo e palestras \u003C3\n\space"souPyLadiesSP \space roadsec",
         "story": "PyLadies S\u00e3o Paulo shared PyLadies Brasil's album.",
        "created time": "2016-11-18T22:22:35+0000",
        "id": "1469656506672137 1615174665453653"
      { ▶ 4 properties, 310 bytes }.
      { ▶ 4 properties, 661 bytes }.
      { > 3 properties, 293 bytes }.
      { > 4 properties, 280 bytes }.
      { > 3 properties, 709 bytes }.
      { ▶ 4 properties, 767 bytes }.
      { > 3 properties, 252 bytes },
      { > 4 properties, 238 bytes }.

↓ 4 properties, 774 bytes }.

      { ▶ 3 properties, 293 bytes }.
      { > 4 properties, 240 bytes }.
      { > 3 properties, 150 bytes }.
      { > 4 properties, 252 bytes }.
      { > 4 properties, 1 KB }.
      { > 3 properties, 483 bytes },
      { > 4 properties, 350 bytes },
      { > 3 properties, 181 bytes },
      { > 4 properties, 448 bytes },
      { > 4 properties, 267 bytes },
      { > 4 properties, 575 bytes },
      { ▶ 3 properties, 252 bytes },
      { ▶ 4 properties, 412 bytes }.
      { > 4 properties, 660 bytes },
      { > 3 properties, 270 bytes }
    "paging": { V 2 properties, 603 bytes
      "previous": "https:\/\graph.facebook.com\/v2.8\/1469656506672137\/posts?since=1479507755&access token=1522873017729918|pum50zN48gJTrwzqZbF5qE9UEc0&limit=25&
      paging token=enc AdAxiPgniWOahSxlvEEosTDbwaamvpHDS0ZBcnUZAhf5BDcz2IuBvernkLeVAfe3b50suZB0WKy80KHnzZAWn0WpNpWVUg8DT2vugfYeSUzoPdFKzgZDZD& previous=1",
      "next": "https:\/\/graph.facebook.com\/v2.8\/1469656506672137\/posts?access token=1522873017729918|pum50zN48gJTrwzqZbF5qE9UEc0&limit=25&until=1474728701&
```

```
# ws_facebook_v5.py (parte 3)
access token = obter access token()
# posts de PyLadiesSP
url = url para obter posts('PyLadiesSP', access token)
condicao = True
while condicao is True:
    dados = ws.coletar dados(url)
    posts = json.loads(dados)
    mensagens = obter mensagens(posts)
    ws.atualizar arquivo('posts.txt', '\n'.join(mensagens))
    url = posts['posts']['paging']['next']
    condicao = ????
???? = vamos pensar quais poderiam ser as condições de parada?
```

Facebook APIs



Tudo tem um id para conseguir obter os dados associados.

A partir do id de um post podemos obter os dados do post, como por exemplo likes:

https://graph.facebook.com/1469656506672137_1615021138802339?field
s=likes&access_token=...

Para saber quais são as várias possibilidades, consulte o Graph API Explorer:

https://developers.facebook.com/tools/explorer

```
Aplicativo: [?]
                                                                                                                                      Graph API Explorer ▼
Graph API Explorer
               6 EAACEdEose0cBAOe3Puzo3FFfQWOwmgVZBwYY7ABaLhZAdgOOIUENo4zkqQc1GmLTmgaLAnYSGvqGtBSxkqVZCDAOHMrfYlk2H9
                                                                                                                                            Token de acesso:
      GET ▼ → /v2.8 ▼ /1469656506672137 1615021138802339?fields=likes
                                                                                                                                                 ▶ Submit
                                                                                                                            Saiba mais sobre a sintaxe da Graph API
  Node: 1469656506672137 16150211
                                         "likes": {
     ✓ likes
                                           "data": [
          + Search for a field
                                               "id": "764628623592259",
                                               "name": "Leila Aparecida da Silva"
      + Search for a field
                                               "id": "840362745983320",
                                               "name": "Caroline Dantas"
                                             },
                                               "id": "770530686319690",
                                               "name": "Eliza Ramos"
                                             },
                                               "id": "761379187251469",
                                               "name": "Débora Fernandes"
                                             },
                                               "id": "1106573342750500",
                                               "name": "Ariadyne Oliveira"
                                             },
                                               "id": "893251710693095",
                                                                                                                                                       51
                                               "name": "Patricia Novais"
                                             },
```





 criar uma função para fazer outras consultas a partir dos posts, por exemplo: likes de cada post.

```
# ws facebook v6.py (parte 3)
def url_para_coletar_likes_mais_recentes_de_um_post(postid):
    return
'https://graph.facebook.com/{postid}?fields=likes&{access token}
 .format(postid=postid, access token=access token)
def coletar likes mais recentes de um post(post):
    url = url para coletar likes mais recentes de um post(
           post['id'])
   resultado likes = coletar dados(url)
   return json.loads(resultado likes)
def obter likes(likes):
   nomes = []
    for like in likes['likes']['data']:
       nomes.append(like['name'])
   return nomes
                                                                   53
```



```
# ws_facebook_v6.py (parte 4)
def post_e_seus_likes(post):
    message = post.get('message', post.get('story', ''))
    likes = coletar_likes_mais_recentes_de_um_post(post)
    return (message, obter_likes(likes))
```

```
# ws facebook v6.py (parte 4)
access token = obter_access_token()
# posts de PyLadiesSP
url = url_para_obter_posts('PyLadiesSP', access_token)
condicao = True
while condicao is True:
    dados = ws.coletar dados(url)
    posts = json.loads(dados)
    for post in posts['posts']['data']:
        msg, likes = post e seus likes(post)
        msg_e_likes = msg_+' \n' + ' \n'.join(likes) + ' \n' + ' - ' * 3 + ' \n'
        ws.atualizar arquivo('posts.txt', msg e likes)
    url = posts['posts']['paging']['next']
    condicao = ????
```

???? = vamos pensar quais poderiam ser as condições de parada?

As meninas do Pyladies dominaram a Vila das Minas na Roadsec esse ano com workshops, dojo e palestras <3

#souPyLadiesSP #roadsec

Naomi Ceder; Daniele Marighetti; Greyce Carla; Fátima Sousa; Patricia Leite; Raquel Barros; Priscila Viana; Eliza Ramos; Juliana Neres; Mazinho Web; Luiz Menezes Filho; Camila Lopes; Patty Vader; Marcelo Yamaki; PyLadies São Paulo; Bruno Rocha;

Já começou o Workshop de Introdução à Python aqui na Roadsec.

E só chegar aqui na área de atividades. ;)

#souPyLadiesSP #roadsec

Leila Aparecida da Silva; Caroline Dantas; Eliza Ramos; Débora Fernandes; Ariadyne Oliveira; Patricia Novais; Mazinho Web; Marcelo Yamaki; Débora Oliveira; Juliana Neres; Natalia Zaramelo de Campos;

Amanhã tem Roadsec e estamos com várias atividades na Vila das Minas. Confira a programação:

10h - Workshop de introdução a Python com as PyLadies São Paulo

11h - Mesa Redonda: Quem são as mulheres na computação com as PyLadies Brasil

13h - Dojo com as PyLadies São Carlos

15h - Workshop de introdução a Machine Learning com as PyLadies Campinas

15h - Palestra PyLadies São Paulo e o que é Python

E ai quem vamos? <3

Vanessa de Moraes; Patricia Novais; Marcelo Yamaki; Leila Aparecida da Silva; Stefanie Souza; Júlia de Almeida; PyLadies Campinas; Mazinho Web; Isis Ramos;

Twitter APIs



- REST (tweets recentes, últimos 7 dias)
- Streaming (tempo real)

\$ terminal: pip install TwitterAPI

Twitter APIs - Preparação



- Criar um arquivo chamado: pizza.py
- Entrar em: https://apps.twitter.com
- Clicar em Create New App
- Criar o app (preencher as informações)

Twitter APIs - Preparação



Create an application

Application Details	
Name *	
PyLadies Test API	
Your application name. This is used to attribute the source of a	tweet and in user-facing authorization screens. 32 characters max.
Description *	
Show Twitter's API	
Your application description, which will be shown in user-facing	authorization screens. Between 10 and 200 characters max.
Website *	
https://beatrizuezu.wordpress.com/	
	ran go to download, make use of, or find out more information about your application. This fully-qualified URL is used in the source attribution
for tweets created by your application and will be shown in user- (If you don't have a URL yet, just put a placeholder here but ren	
Caliback URL	
Where should we return after successfully authenticating? OAut your application from using callbacks, leave this field blank.	h 1.0a applications should explicitly specify their oauth_callback URL on the request token step, regardless of the value given here. To restrict
you approved their same wantering, 1987 C the Heat Meth.	

Developer Agreement

Yes, I have read and agree to the Twitter Developer Agreement.

Twitter APIs - Preparação



- Clicar em 'Manage keys and access tokens'
- Pegar o consumer_key e consumer_secret
- Mais pra baixo da página terá um botão para gerar o Token, gerar tokens
- Pegar o access_token e o access_token_secret
- No arquivo pizza.py coloque as variáveis acima com seus respectivos valores.

Twitter APIs - REST



from TwitterAPI import TwitterAPI, TwitterRestPager

```
search term = 'pizza'
consumer key = # sua consumer key
consumer secret = # sua consumer secret
access token = # sua access token
access token secret = # sua access token secret
api = TwitterAPI(consumer key,consumer secret, access token,
access token secret)
resultado = TwitterRestPager(api, 'search/tweets', {'q': search term})
for item in resultado.get iterator():
   r = item['text'] if 'text' in item else item
    print(item['created at'])
   print(r.encode(encoding=sys.getdefautencoding()))
```

Twitter APIs - Streaming



• Crie um outro arquivo chamado: mais_pizza.py

Twitter APIs - Streaming



```
from TwitterAPI import TwitterAPI

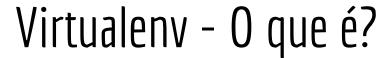
search_term = 'pizza'
consumer_key = # sua consumer key
consumer_secret = # sua consumer secret
access_token = # sua access token
access_token_secret = # sua access token secret

api = TwitterAPI(consumer_key,consumer_secret, access_token, access_token_secret)
resultado = api.request('statuses/filter', {'track': search_term})
```



```
for item in resultado:
    try:
        print(item['created_at'])
        print(item['user']['location'])
        print(item['user']['name'])
        print(r)
        print('-')
    except:
        pass
```

https://gist.githubusercontent.com/hrp/900964/raw/2bbee4c296e6b54877b537144be89f19beff75f4/twitter.json





Virtualenv é uma ferramenta para criar ambientes Python isolados.

Virtualenv cria uma pasta que contém todos os executáveis necessários para um projeto Python.





- 1. Cria uma pasta para seu projeto
- 2. Entre na pasta
- 3. Crie o ambiente virtual

python -m venv .venv

4. Ative o ambiente virtual

.venv\scripts\activate.bat