

**TECNOLOGIA EM SISTEMAS PARA INTERNET**

**Daniel Evangelista Pereira  
Ribson Coelho Cardoch Valdés  
Douglas Seidi Shibata**

**RELATÓRIO DE PRÁTICA INTEGRADA  
DE  
CIÊNCIA DE DADOS E INTELIGÊNCIA ARTIFICIAL**

19/08/2020

**Brasília - DF**

**2020**

# Sumário

<b>1. Objetivos</b>	<b>3</b>
<b>2. Descrição do problema</b>	<b>4</b>
<b>3. Desenvolvimento</b>	<b>5</b>
3.1 Código implementado	5
<b>4. Considerações Finais</b>	<b>10</b>
<b>Referências</b>	<b>11</b>

# 1. Objetivos

Esta etapa do projeto tem como finalidade a exploração dos dados utilizando a linguagem SQL no script do notebook em python, que foram coletados na fase anterior por meio de web scraping, utilizando a linguagem python e gerando um arquivo Comma-separated values(CSV). Que será utilizado para a exploração. E o código estará armazenado no repositório do Github com a utilização do versionamento de código do git.

## 2. Descrição do problema

Ao iniciar o desenvolvimento do projeto podemos destacar a forma da coleta dos dados e como estão distribuídos no site e a partir desta análise inicial, foi visto que havia uma grande quantidade de dados em diversas páginas que precisar ser extraídas do site para realizar a coleta de dados e armazenar em DataFrames para que se possa gerar um arquivo csv para começar o estudo dos dados na primeira etapa do projeto. E a partir dos dados que foram gerados e armazenados no csv será feita a exploração de dados com a utilização da linguagem sql para a melhor filtragem dos dados.

## 3. Desenvolvimento

As tecnologias utilizadas para a elaboração da primeira parte do projeto na fase de coleta de dados consiste na linguagem de programação python com o auxílio de suas bibliotecas, como o requests para fazer a requisição do site e retornar o status code, a biblioteca BeautifulSoup para a leitura das páginas em Hyper Text Markup Language(HTML) e para o armazenamento de dados foi utilizada a biblioteca pandas, e o pandasql para que seja possível executar códigos SQL para fazer consultas no dataframe de ovnis, por meio do notebook e o ambiente de desenvolvimento Google Colab.

### 3.1 Código implementado

#### 5.4 - Exploração de dados com SQL

##### - 5.4 - Exploração dos dados com SQL

1. Saber a quantidade de linhas, observações ou variáveis que foram coletadas.
2. Quantos relatos ocorreram por estado em ordem decrescente?
3. Remover possíveis campos vazios (sem estado).
4. Limitar a análise aos estados dos Estados Unidos.
5. Consulta por cidades, com o objetivo de saber quais contêm o maior número de relatos (cidades que apresentem ao menos 10 relatos).
6. Com o dado anterior, responder a seguinte pergunta: por que será que essa é a cidade que possui mais relatos?
7. Fazer uma query exclusiva para o estado com maior número de relatos, buscando cidades que possuam um número superior a 10 relatórios. Enfatizar a cidade, a quantidade de relatos e formato do objeto não identificado.

```
[29] 1 #Para baixar o pandas sql, apenas tirar a # e executar a célula  
     2 !pip install pandasql
```

```
Requirement already satisfied: pandasql in /usr/local/lib/python3.6/dist-packages (0.7.3)  
Requirement already satisfied: sqlalchemy in /usr/local/lib/python3.6/dist-packages (from pandasql) (1.3.19)  
Requirement already satisfied: numpy in /usr/local/lib/python3.6/dist-packages (from pandasql) (1.18.5)  
Requirement already satisfied: pandas in /usr/local/lib/python3.6/dist-packages (from pandasql) (1.0.5)  
Requirement already satisfied: pytz>=2017.2 in /usr/local/lib/python3.6/dist-packages (from pandas->pandasql) (2018.9)  
Requirement already satisfied: python-dateutil>=2.6.1 in /usr/local/lib/python3.6/dist-packages (from pandas->pandasql) (2.8.1)  
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.6/dist-packages (from python-dateutil>=2.6.1->pandas->pandasql) (1.15.0)
```

Baixando a biblioteca pandasql para a utilização de códigos SQL na exploração de dados

```
1 #Importando a biblioteca pandasql e pandas  
2 import pandasql  
3 import pandas as pd
```

Importação das bibliotecas pandasql e pandas para exploração de dados com sql e a leitura e geração de arquivos csv, respectivamente.

```
[25] 1  ovnis = pd.read_csv('OVNIS.csv')
      2  ovnis
```

	Date / Time	City	State	Shape	Duration	Summary	Posted
0	9/30/97 20:00	Nova Scotia (Canada)	NS	Light	8-10 seconds.	Single light resembling a star, but moving spu...	10/30/06
1	9/22/97 20:00	Solomons Island	MD	Disk	10 minutes	Close up at twilight, Stationary UFO.	8/5/09
2	9/19/97	Garden Grove	CA	Rectangle	4 mins.	Around 6:30 PM I was walking through a Vons Pa...	12/1/19
3	9/18/97 20:15	Panama City	FL	Unknown	30 seconds	Looked like stars in the sky so far up/moveing...	3/13/12
4	9/15/97 00:00	Houston	TX	Disk	5 minutes	Beautiful silver-colored flying saucer about t...	7/19/10
...	...	...	...	...	...	...	...
71896	8/1/17 02:45	Corcoran	MN	Light	Still going	Small light south west of Minneapolis maneuver...	8/4/17
71897	8/1/17 02:00	Moreno Valley	CA	Other	10 seconds	I was looking out the front windshield and loo...	8/4/17
71898	8/1/17 01:00	Bradenton	FL	Other	<20 seconds	I was walking my dog about 1am on August 1, 20...	5/9/19
71899	8/1/17 00:00	Springdale	AR	NaN	1 hour	Glowing flying people . seven of them flying I...	2/13/20
71900	8/1/17	Laurel	MD	Other	NaN	It was an alien project level 1 federal ran on...	6/25/20

71901 rows x 7 columns

Leitura do arquivo OVNIS csv e armazenando na variável ovnis

```
[26] 1  #Saber a quantidade de linhas, observações ou variáveis que foram coletadas.
      2  print("Quantidade de Linhas coletadas: ",len(ovnis))
```

```
Quantidade de Linhas coletadas: 71901
```

Imprimindo a quantidade de linhas que foram coletadas na atividade anterior

```
1  #Quantos relatos ocorreram por estado em ordem decrescente?
2  estados =ovnis.State.value_counts()
3  estados.sort_values(ascending=False)
```

```
CA    7911
FL    4352
WA    3225
TX    2882
NY    2824
...
NF      21
YT      14
PE       9
NT       7
SA       4
Name: State, Length: 64, dtype: int64
```

Declarando a variável estados e contando os estados por valores, ou seja, por siglas que apresentam na base de dados e depois organizando as siglas por ordem crescente.

```

1 #Remover possíveis campos vazios (sem estado).
2 ovis.State.dropna()
3 ovis.State.sort_values(ascending=True).unique()

array(['AB', 'AK', 'AL', 'AR', 'AZ', 'BC', 'CA', 'CO', 'CT', 'DC', 'DE',
      'FL', 'GA', 'HI', 'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA', 'MA',
      'MB', 'MD', 'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NB', 'NC', 'ND',
      'NE', 'NF', 'NH', 'NJ', 'NM', 'NS', 'NT', 'NV', 'NY', 'OH', 'OK',
      'ON', 'OR', 'PA', 'PE', 'QC', 'RI', 'SA', 'SC', 'SD', 'SK', 'TN',
      'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV', 'WY', 'YT', nan],
      dtype=object)

```

Removendo os possíveis campos vazios da base de dados.

```

[31] 1 #Limitar a análise aos estados dos Estados Unidos.
2 q = ""
3 | SELECT * from total where State LIKE '%AK%' OR State LIKE '%AL%'
4 | OR State LIKE '%AR%'
5 | OR State LIKE '%AZ%'
6 | OR State LIKE '%CA%'
7 | OR State LIKE '%CO%'
8 | OR State LIKE '%CT%'
9 | OR State LIKE '%DE%'
10 | OR State LIKE '%FL%'
11 | OR State LIKE '%GA%'
12 | OR State LIKE '%HI%'
13 | OR State LIKE '%IA%'
14 | OR State LIKE '%ID%'
15 | OR State LIKE '%IL%'
16 | OR State LIKE '%IN%'
17 | OR State LIKE '%KS%'
18 | OR State LIKE '%KY%'
19 | OR State LIKE '%LA%'
20 | OR State LIKE '%MA%'
21 | OR State LIKE '%MD%'
22 | OR State LIKE '%ME%'
23 | OR State LIKE '%MI%'
24 | OR State LIKE '%MN%'
25 | OR State LIKE '%MO%'
26 | OR State LIKE '%MS%'
27 | OR State LIKE '%MT%'
28 | OR State LIKE '%NC%'
29 | OR State LIKE '%ND%'
30 | OR State LIKE '%NE%'
31 | OR State LIKE '%NH%'
32 | OR State LIKE '%NJ%'
33 | OR State LIKE '%NM%'
34 | OR State LIKE '%NV%'
35 | OR State LIKE '%NY%'
36 | OR State LIKE '%OH%'
37 | OR State LIKE '%OK%'
38 | OR State LIKE '%OR%'
39 | OR State LIKE '%PA%'
40 | OR State LIKE '%RI%'
41 | OR State LIKE '%SC%'
42 | OR State LIKE '%SD%'
43 | OR State LIKE '%TN%'
44 | OR State LIKE '%TX%'
45 | OR State LIKE '%UT%'
46 | OR State LIKE '%VA%'
47 | OR State LIKE '%VT%'
48 | OR State LIKE '%WA%'
49 | OR State LIKE '%WI%'
50 | OR State LIKE '%WV%'
51 | OR State LIKE '%WY%'
52 | ""
53 # Executa o seu comando SQL e retorna um dataframe
54 just_us = pandasql.sqldf(q.lower(), locals())
55 just_us

```

Limitando os estados para apenas estados norte americanos e executando o comando sql.



	Date / Time	City	State	Shape	Duration	Summary	Posted
0	9/22/97 20:00	Solomons Island	MD	Disk	10 minutes	Close up at twilight, Stationary UFO.	8/5/09
1	9/19/97	Garden Grove	CA	Rectangle	4 mins.	Around 6:30 PM I was walking through a Vons Pa...	12/1/19
2	9/18/97 20:15	Panama City	FL	Unknown	30 seconds	Looked like stars in the sky so far up/moveing...	3/13/12
3	9/15/97 00:00	Houston	TX	Disk	5 minutes	Beautiful silver-colored flying saucer about t...	7/19/10
4	9/15/97 20:00	Santa Fe	NM	Light	2-3 minutes	Saw white dot of light moving in zig-zag motio...	11/9/17
...	...	...	...	...	...	...	...
64892	8/1/17 02:45	Corcoran	MN	Light	Still going	Small light south west of Minneapolis maneuver...	8/4/17
64893	8/1/17 02:00	Moreno Valley	CA	Other	10 seconds	I was looking out the front windshield and loo...	8/4/17
64894	8/1/17 01:00	Bradenton	FL	Other	<20 seconds	I was walking my dog about 1am on August 1, 20...	5/9/19
64895	8/1/17 00:00	Springdale	AR	None	1 hour	Glowing flying people . seven of them flying l...	2/13/20
64896	8/1/17	Laurel	MD	Other	None	It was an alien project level 1 federal ran on...	6/25/20

64897 rows x 7 columns

Resultado da consulta feita para filtrar apenas estados norte americanos.

```
[35] 1 just_us.to_csv('dados_usa.csv')
```

Gerando um arquivo csv dos dados que foram gerados anteriormente, com a filtragem de apenas estados americanos.

```
[ ] 1 cidades_usa = just_us.City.value_counts()
    2 cidades_usa[cidades_usa >=10]
    3 cidades_usa.head()
```



```
Phoenix      366
Las Vegas    338
Seattle      323
Portland     318
San Diego    272
Name: City, dtype: int64
```

Contando a quantidade de caso por cidades dos estados unidos, e mostrando apenas os estados que têm mais de 10 ocorrências. E mostrando apenas os 5 mais recorrentes.



```
1 #Fazer uma query exclusiva para o estado com maior número de relatos, buscando cidades que possuam um número superior a 10 relatos.
2 #Enfatizar a cidade, a quantidade de relatos e formato do objeto não identificado.
3
4 # Executa o seu comando SQL e retorna um dataframe
5 query = '''
6 SELECT CITY,Count(City) as 'Numero de Ocorrencias', Shape FROM just_us group by city,shape having count(City)>=10 order by count(City) desc
7 '''
8 pandasql.sqldf(query.lower(), locals())
9
```

	City	numero de ocorrencias	Shape
0	Seattle	75	Light
1	Las Vegas	74	Light
2	Phoenix	72	Light
3	San Diego	71	Light
4	Myrtle Beach	66	Light
...	...	...	...
564	Virginia Beach	10	Sphere
565	Visalia	10	Light
566	Wellington	10	Light
567	Woodbridge	10	Light
568	York	10	Circle

569 rows x 3 columns

Identificando os formatos dos objetos não identificados por cidade e quantidade de ocorrência do formato que foi relatado.

## 4. Considerações Finais

Nesta parte de exploração de dados, houve estudos das bibliotecas utilizadas e o modo que as consultas foram realizadas por meio da linguagem de consulta estruturada. No processo de desenvolvimento houve alguns empecilhos, como a forma de executar alguns comandos SQL para que se obtenha o exato resultado e a remoção de dados nulos, houve uma pequena dificuldade, pois acredita-se que alguns campos não estejam nulos, mas está escrito com sem valor.

## Referências

ANDRADE, Andrew. Análise Exploratória de Dados.**Escola de Dados**. [s.d.]. Disponível em: <<https://escoladedados.org/tutoriais/analise-exploratoria-de-dados/>>. Acesso em: 13, Setembro de 2020.

GEROLA, Letícia. 5 pythons hacks para a exploração do dados.**Medium**, 2020. Disponível em: <<https://medium.com/joguei-os-dados/4-python-hacks-para-explora%C3%A7%C3%A3o-de-dados-8c89931b6d1f>>. Acesso em: 13, Setembro de 2020.

CARACIOLO, Marcel. **Introdução a análise exploratória com python e pandas**. 2013.(1h32m25s). Disponível em: <<https://www.youtube.com/watch?v=vlJwq6QjZL8>>. Acesso em: 14, Setembro de 2020.