

## Relatório UFO Data – Parte 2

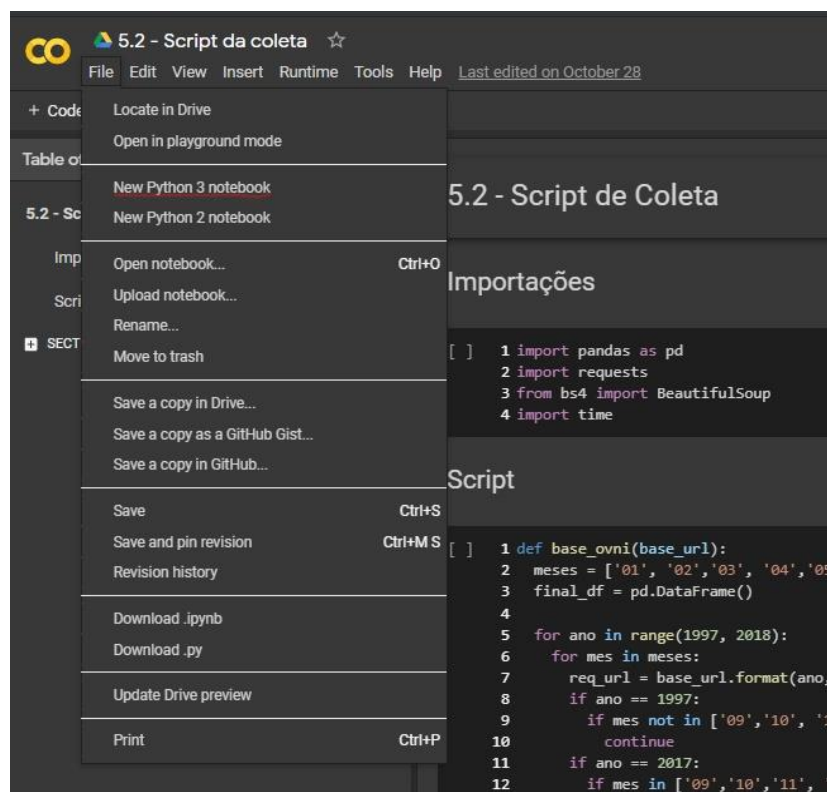
Levi Alves de Freitas Junior

### Problema

Criar uma análise utilizando banco de dados, aproveitando a base de dados requisitada no primeiro relatório, e utilizar o pandasql para a exploração dos dados. Vamos utilizar a base gerada no primeiro relatório e explora-la.

Para a primeira parte vamos criar um novo notebook no colab research. Siga as instruções a seguir e veja as imagens ilustrativas:

**1** – Logo após terminar as instruções da primeira etapa de análise dos relatos dos ovni's, crie um novo notebook no colab, acessando no menu a opção File e logo após vá em New Python 3 notebook.



**2** – Após criar um novo notebook, podemos começar a criar nosso código e inicialmente precisamos importar as bibliotecas pandas e para usar o banco de dados vamos importar o pandasql.

```
#!pip install -U pandasql
import pandas as pd
import pandasql
```

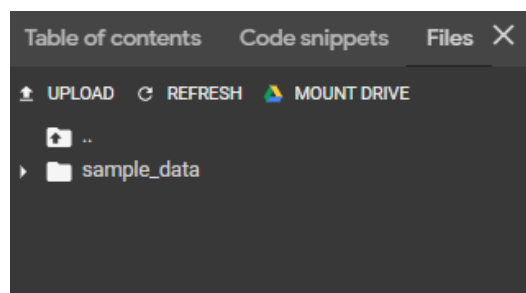
Primeiro selecione o comando - **!pip install -U pandasql** e aperte CTRL+ SHIFT + ENTER para instalar o pandasql no colab, logo após inicialize a célula no colab para importar o pandas e o pandasql,

Obs. O import do pandasql pode não funcionar antes de utilizar o comando pip para utilizá-lo.

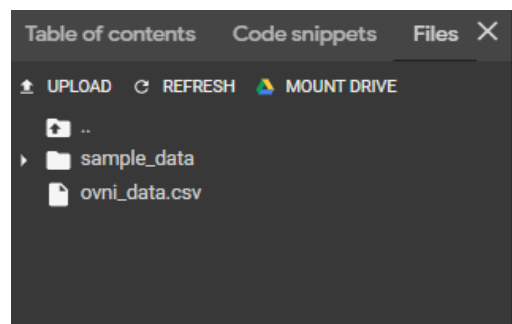
### 3 – Agora vamos importar a base de dados

Lembre-se de importar a base de dados no colab para podermos chama-la no código, siga as instruções abaixo para colocar a base.

No menu esquerdo selecione Files, e clique em upload, selecione a base, onde você salvou



Aguarde um pouco e sua base estará inserida no colab.



Podemos prosseguir agora e chamar nossa base de dados.

```
df = pd.read_csv('ovni_data.csv')

#Excluindo coluna Indesejada
del df['Unnamed: 0']

#df.columns = ['Date/Time','City','State','Shape','Duration','Summary','Posted']

linhas = df['Date / Time'].size

print("Linhas: ",linhas)
```

Com o comando 'del' retiramos a coluna não nomeada pois não será necessária, então vamos imprimir a quantidade de linhas do nosso dataframe.

O resultado deve ser:

Linhas: 99705

#### 4 – Vamos organizar os dados e remover campos em branco

```
estados = """
    SELECT State, COUNT(State) as qtd_rel FROM df WHERE State != ' ' GROUP BY State ORDER BY COUNT(State) desc ;
    """
result_consulta = pandasql.sqldf(estados)
result_consulta
```

Fazemos uma pesquisa utilizando o pandasql, os estados com mais relatos de aparições de ovni's e organizamos em ordem decrescente.

O estado com mais relato é CA – Califórnia com 11403 relatos.

	State	qtd_rel
0	CA	11403
1	FL	5577
2	WA	4901
3	TX	4154
4	NY	3871
...	...	...
63	PE	20
64	YT	19
65	PR	18
66	YK	5
67	VI	1

68 rows × 2 columns

5 – Neste momento vamos analisar somente os casos relatados nos Estados Unidos, com o seguinte código vamos comparar com a base de dados para filtrar e retornar os relatos.

```
#Cria um array com os estados da coluna estado.

lista_usa = ['AK', 'AL', 'AR', 'AZ', 'CA', 'CO', 'CT', 'DE', 'FL', 'GA',
             'HI', 'IA', 'ID', 'IL', 'IN', 'KS', 'KY', 'LA', 'MA', 'MD',
             'ME', 'MI', 'MN', 'MO', 'MS', 'MT', 'NC', 'ND', 'NE', 'NH',
             'NJ', 'NM', 'NV', 'NY', 'OH', 'OK', 'OR', 'PA', 'RI', 'SC',
             'SD', 'TN', 'TX', 'UT', 'VA', 'VT', 'WA', 'WI', 'WV', 'WY']

lista = pd.DataFrame(lista_usa, columns=['Estados_Unidos'])
```

Criamos uma lista com as iniciais do estado e armazenamos esta lista à variável lista\_usa e logo após criamos uma coluna de um dataframe e armazenamos na variável lista.

Com a coluna 'Estados\_Unidos' do dataframe criada podemos agora comparar o dataframe original – 'df' com a coluna gerada e retornar somente os dados dos estados que constam na lista\_usa. Então fazemos a seguinte query.

```
consulta_est_validos = '''
    SELECT df.* FROM df, lista
    WHERE State = Estados_Unidos
    '''

limit_usa = pandasql.sqldf(consulta_est_validos)
```

O resultado será:

	Date / Time	City	State	Shape	Duration	Summary	Posted
0	9/30/97 22:00	Madison	WI	Light	5 minutes	Strange light inside Lake Monona	3/2/04
1	9/28/97 23:15	San Francisco	CA	Triangle	12-15s	flying-wing shape outlined by 12-14 lights. Ap...	7/5/99
2	9/27/97 23:00	Egan	SD	Other	30 minutes	The Weirdest Thing I Have Ever Seen	2/22/05
3	9/27/97 05:00	Crestwood	KY	Disk	15 minutes	A big disk with red and green lights on the ri...	8/5/01
4	9/25/97 22:00	Clearfield	UT	Triangle	60-90 seconds	We observed a low flying craft (aprox. 100yards...	1/28/99
...	...	...	...	...	...	...	...
88137	8/1/17 14:00	Joliet	IL	Other	2 minutes	The White Cube UFO	7/25/19
88138	8/1/17 06:15	Columbus (North)	GA	Fireball	3 seconds	Green streak growing in size moving from west ...	8/4/17
88139	8/1/17 02:45	Corcoran	MN	Light	Still going	Small light south west of Minneapolis maneuver...	8/4/17
88140	8/1/17 02:00	Moreno Valley	CA	Other	10 seconds	I was looking out the front windshield and loo...	8/4/17
88141	8/1/17 01:00	Bradenton	FL	Other	<20 seconds	I was walking my dog about 1am on August 1, 20...	5/9/19

## 6 – Consulta e filtragem das cidades com 10 relatos ou mais

```

cidades_agrup = """
    SELECT State, Shape, City, COUNT(*) as total_posts FROM limit_usa GROUP BY City;
    """
result_cidades_agrup = pandasql.sqldf(cidades_agrup)

nova_query_usa = """
    SELECT State, Shape, City, total_posts
    FROM result_cidades_agrup
    WHERE City != 'unknown' AND total_posts >= 10
    ORDER BY total_posts desc;
    """
sql_cidades_usa = pandasql.sqldf(nova_query_usa)

```

Utilizamos a primeira query para agrupar as cidades e criamos uma subquery para tirar os campos não nomeados, no caso, o comando AND, para que o total de posts (relatos) seja igual ou maior que 10.

O resultado desta query deve ser como no modelo abaixo.

	State	Shape	City	total_posts
0	AZ	Flash	Phoenix	558
1	WA	Changing	Seattle	548
2	OR	Circle	Portland	480
3	NV	Rectangle	Las Vegas	473
4	CA	Rectangle	San Diego	394
...	...	...	...	...
1833	AZ	Changing	Winslow	10
1834	MT	Flash	Wolf Point	10
1835	TX	Circle	Woodville	10
1836	MA	Teardrop	Yarmouth	10
1837	CA	Disk	Yucca Valley	10
1838 rows × 4 columns				

7 – Porque esta é a cidade que mais tem relatos.

**Phoenix** - No dia 13 de março de 1997, o céu do estado do Arizona e de Nevada, encheram de luzes, milhares de pessoas presenciaram o ocorrido, as pessoas afirmaram ver algo em formato triangular vagando pelo céu na cidade de Phoenix, capital do Arizona e cidade com mais relatos na análise de todos os casos de aparecimento de OVNI's, o próprio governador presenciou e afirmou que seria um "objeto de outro mundo". Milhares de pessoas registraram essas aparições isso explica ser um dos locais com mais relatos.

8 – Busca nos estados com maior número de relato, juntamente buscando as cidades que tenham um número de relatos superior a 10 relatórios.

```
estados_usa_mais_rel = """
    SELECT State, MAX(qtd_rel) as qtd_rel FROM result_consulta;
    """

sql_estados = pandasql.sqldf(estados_usa_mais_rel)

cid_mais_rel = """
    SELECT sql_cidades_usa.City, sql_cidades_usa.total_posts, sql_cidades_usa.Shape
    FROM sql_estados, sql_cidades_usa
    WHERE sql_cidades_usa.State = sql_estados.State AND sql_cidades_usa.total_posts > 10 AND sql_cidades_usa.Shape NOT IN (' ', 'Unknown')
    """

result_final_cidades = pandasql.sqldf(cid_mais_rel)
```

O resultado deve ser o seguinte rodando result\_final\_cidades.

	City	total_posts	Shape
0	San Diego	394	Rectangle
1	Los Angeles	379	Changing
2	Sacramento	242	Changing
3	San Jose	223	Other
4	San Francisco	200	Chevron
...	...	...	...
221	Imperial Beach	11	Fireball
222	Miramar	11	Formation
223	Nipomo	11	Fireball
224	Pismo Beach	11	Disk
225	West Los Angeles	11	Triangle
226 rows × 3 columns			

Este resultado é em cima de uma query onde são retirados campos em branco e campos Unknown. Abaixo faremos uma consulta com os dois campos o que irá resultar em uma pesquisa e um resultado diferente.

```

estados_usa_mais_rel = """
    SELECT State, MAX(qtd_rel) as qtd_rel FROM result_consulta;
    """

sql_estados = pandasql.sqldf(estados_usa_mais_rel)

cid_mais_rel = """
    SELECT sql_cidades_usa.City, sql_cidades_usa.total_posts, sql_cidades_usa.Shape
    FROM sql_estados, sql_cidades_usa
    WHERE sql_cidades_usa.State = sql_estados.State AND sql_cidades_usa.total_posts > 10
    """

result_final_cidades = pandasql.sqldf(cid_mais_rel)

```

Somente retiramos o final do código da query, retiramos tudo a partir do último AND do WHERE até o fim da query.

Como resultado deveremos receber o seguinte dataframe:

	City	total_posts	Shape
0	San Diego	394	Rectangle
1	Los Angeles	379	Changing
2	Sacramento	242	Changing
3	San Jose	223	Other
4	San Francisco	200	Chevron
...	...	...	...
240	Nipomo	11	Fireball
241	Pismo Beach	11	Disk
242	Tehachapi	11	Unknown
243	West Los Angeles	11	Triangle
244	Wildomar	11	Unknown
245 rows × 3 columns			

## Referências

Relatos dos ovni's - <https://aventurasnahistoria.uol.com.br/noticias/reportagem/luzes-phoenix-bizarras-aparicoes-de-ovnis-nos-estados-unidos.phtml>

O projeto encontra-se no GitHub – para eventual consulta.

[https://github.com/LeviAFJunior/Analise\\_Dados\\_Base\\_ovni/tree/master/Base\\_Ovni\\_parte2](https://github.com/LeviAFJunior/Analise_Dados_Base_ovni/tree/master/Base_Ovni_parte2)