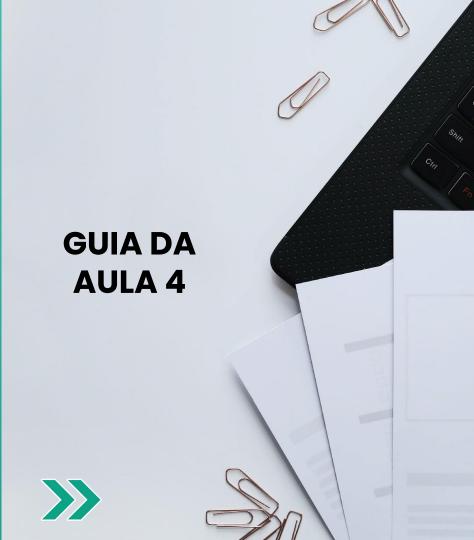


# Profissão: Analista de dados





# PROJETO: ANÁLISE EXPLORATÓRIA DE DADOS EM PYTHON







# Transformação e limpeza de dados

- Dados
- Enriquecimento

Qualidade



Acompanhe aqui os temas que serão tratados na videoaula







#### 1. Dados

O dado bruto é um arquivo do tipo JSON com uma lista de instâncias de entregas. Cada instância representa um conjunto de **entregas** que devem ser realizadas pelos **veículos** do **hub** regional. Exemplo:

```
"igson [ { "name": "cvrp-0-df-0", "region": "df-0", "origin": {"Ing": -47.802664728268745, "lat": -15.657013854445248}, "vehicle_capacity": 180, "deliveries": [ { "id": "ed0993f8cc70d998342f38ee827176dc", "point": {"Ing": -47.7496622016347, "lat": -15.65879313293694}, "size": 10 }, { "id": "c7220154adc7a3def8f0b2b8a42677a9", "point": {"Ing": -47.75887552060412, "lat": -15.651440380492554}, "size": 10 }, ... ] } ] ...
```





Processamos o dado bruto e construímos o DataFrame Pandas deliveries \_df através de operações como achatamento (flatten) e explosão (explode) de colunas:

```
In []:
    !wget -q << EOF
    https://raw.githubusercontent.com/andre-marcos-perez/ebac-course-
    utils/main/dataset/deliveries.json
    EOF \
    -O deliveries.json

In []:
    import json
    import pandas as pd
    # dado bruto em um dict</pre>
```





```
with open ('deliveries.json', mode='r', encoding='utf8') as file:
 data = json.load(file)
# dado bruto no pandas
deliveries df = pd.DataFrame (data)
# coluna origin
hub origin df = pd.json normalize (deliveries df ["origin"])
deliveries df = pd.merge(left=deliveries df, right=hub origin df,
                         how='inner', left index=True, right index=True)
deliveries df = deliveries df.drop("origin", axis=1)
deliveries df = deliveries df[
    ["name", "region", "lng", "lat", "vehicle capacity", "deliveries"]
deliveries df.rename(columns={"lng": "hub lng", "lat": "hub lat"},
                     inplace=True)
```





```
# coluna deliveries
deliveries exploded df = deliveries df [["deliveries"]].explode("deliveries")
deliveries normalized df = pd.concat([
 pd.DataFrame (deliveries exploded df ["deliveries"].apply(
      lambda record: record["size"]
  )).rename(columns={"deliveries": "delivery size"}),
  pd.DataFrame (deliveries exploded df ["deliveries"].apply(
      lambda record: record["point"]["lng"]
  )).rename(columns={"deliveries": "delivery lng"}),
 pd.DataFrame (deliveries exploded df ["deliveries"].apply(
      lambda record: record["point"]["lat"]
 )).rename(columns={"deliveries": "delivery lat"}),
1, axis=1)
deliveries df = deliveries df.drop("deliveries", axis=1)
deliveries df = pd.merge(left=deliveries df, right=deliveries normalized df,
                         how='right', left index=True, right index=True)
deliveries df.reset index (inplace=True, drop=True)
```

```
In [ ]: deliveries_df.head()
```





## 2. Enriquecimento

## Geocodificação reversa do hub

A geocodificação é o processo que transforma uma localização descrita por um texto (endereço, nome do local etc.) em sua respectiva coodernada geográfica (latitude e longitude). A geocodificação reversa faz o oposto, transforma uma coordenada geográfica de um local em suas respectivas descrições textuais.





Empresas como Google, Bing e Yahoo! fornecem geocodificação como serviço (e cobram por isso). Existe uma projeto *open source* chamado de <u>OpenStreetMap</u> que mantem um serviço gratuito de geocodificação chamado <u>Nominatim</u>, serviço este que apresenta como limitação a quantia de <u>uma única consuta por segundo</u>. Vamos utilizá-lo através do pacote Python geopy para fazer a operação reversa e enriquecer o nosso DataFrame principal.





```
import json
import geopy
from geopy.geocoders import Nominatim

geolocator = Nominatim(user_agent="ebac_geocoder")
location = geolocator.reverse("-15.657013854445248, -47.802664728268745")

print(json.dumps(location.raw, indent=2, ensure_ascii=False))
```

Vamos, então, aplicar a geocodificação nas coordenadas das três regiões e extrair informações de **cidade** e **bairro**.

```
In []:
    from geopy.extra.rate_limiter import RateLimiter

geocoder = RateLimiter (geolocator.reverse, min_delay_seconds =1)
```





```
In [ ]:
         hub df["coordinates"] = hub df["hub lat"].astype(str) + ", "
         hub df["coordinates"] = hub df["coordinates"] + hub df["hub lng"].astype(str)
         hub df["geodata"] = hub df["coordinates"].apply(geocoder)
         hub df.head()
In [ ]:
         hub geodata df = pd.json normalize (hub df ["geodata"].apply (
             lambda data: data.raw
         ) )
         hub geodata df.head()
In [ ]:
         import numpy as np
         hub geodata df = hub geodata df [
              ["address.town", "address.suburb", "address.city"]
```





```
hub geodata df .rename (
    columns={
        "address.town": "hub town",
        "address.suburb": "hub suburb",
        "address.city": "hub city"
    }, inplace=True)
hub geodata df ["hub city"] = np.where(
    hub geodata df ["hub city"].notna(),
    hub geodata df ["hub city"],
    hub geodata df ["hub town"]
hub geodata df ["hub suburb"] = np.where(
    hub geodata df ["hub suburb"].notna(),
   hub geodata df ["hub suburb"],
    hub geodata df ["hub city"])
hub geodata df = hub geodata_df.drop("hub_town", axis=1)
hub geodata df .head()
```





O DataFrame hub\_geodata\_df com as informações de cidade e bairro é, então, combinado ao DataFrame principal deliveries df enriquecendo assim o dado.

```
In []: hub_df = pd.merge(
        left=hub_df, right=hub_geodata_df, left_index=True, right_index=True)
        hub_df = hub_df[["region", "hub_suburb", "hub_city"]]
        hub_df.head()
```





```
In [ ]:
         deliveries df = pd.merge(
             left=deliveries df, right=hub df, how="inner", on="region"
         deliveries df = deliveries df[[
             "name",
             "region",
             "hub lng",
             "hub lat",
             "hub city",
             "hub suburb",
             "vehicle_capacity",
             "delivery size",
             "delivery lng",
             "delivery lat"
         deliveries df.head()
```





#### · Geocodificação reversa da entrega

Enquanto o **hub** contem apenas **3** geolocalizações distintas, as **entregas** somam o total de **636.149**, o que levaria em torno de 7 dias para serem consultadas no servidor do Nominatim, dada a restrição de uma consulta por segundo. Contudo, para cargas pesadas como esta, o *software* oferece uma instalação local (na sua própria máquina) que pode ser utilizada sem restrição.

Confira o link de instalação local:

https://nominatim.org/release-docs/latest/admin/Installation/





Como a instalação do servidor local envolve tecnologias que estão fora do escopo deste curso, como <u>Docker</u>, vamos providenciar estes dados através do *link*:

https://raw.githubusercontent.com/andre-marcos-perez/ebac-course-utils/main/dataset/deliveries-geodata.csv

```
In []:
    !wget -q << EOF
    https://raw.githubusercontent.com/andre-marcos-perez/ebac-
    course-utils/main/dataset/deliveries-geodata.csv
    EOF \
    -O deliveries-geodata.csv</pre>
```









### 3. Qualidade

Qualidade do dados está relacionado a consistência do seu schema valores faltantes etc.

```
In []: deliveries_df.info()

In []: deliveries_df.isna().any()
```





#### · Geocodificação reversa

```
In [ ]:
         100 * (
         deliveries df ["delivery city"].isna().sum() / len(deliveries df)
In [ ]:
         100 * (
         deliveries df ["delivery suburb"].isna().sum() / len(deliveries df)
In [ ]:
         prop df = deliveries df [["delivery city"]].value counts (
         ) / len(deliveries df)
         prop df.sort values (ascending=False).head(10)
In [ ]:
         prop df = deliveries df [["delivery suburb"]].value counts (
         ) / len(deliveries df)
         prop df.sort values (ascending=False) .head(10)
```

