Cálculo de distância entre duas geolocalizações

None

Table of contents

1. Cá	álculo de distância entre duas geolocalizações	3
1.1	Estrutura do projeto	3
1.2	Para o projeto	3
1.3	Links úteis	3
2. Sc	cript - Limpeza e formatação de CEPs	4
2.1	Importando dados	4
2.2	Aplicando tratamento nos dados de CEP	4
2.3	Exportando CEPs limpos	4
3. Sc	cript - Busca de geolocalização	5
3.1	Importando dados limpos	5
3.2	Criando array apenas com CEPs distintos	5
3.3	Buscando geolocalização	5
4. Sc	cript - Calcular distância das geolocalizações	6
4.1	Importando dados	6
4.2	Dando merge com os CEPs conhecidos	6
4.3	Cálculo de distância	6
5. Sc	cript - Funções úteis	7
5.1	FaixaCep	7
5.2	GeoCode	8
5.3	calcular_distancia(origin, destination, metric='hav')	9
5.4	formatar_cep(cep, faixacep)	10
5.5	limpar_cep(cep)	11

- 2/12 - ScientifiCloud

1. Cálculo de distância entre duas geolocalizações

O projeto foi feito para calcular a distância entre duas geolocalizações.

As etapas são construídas desde a limpeza e formatação de CEPs, utilização da API Google para buscar as geolocalizações dos CEPs e o cálculo da distância entre as geolocalizações obtidas.

1.1 Estrutura do projeto

```
mkdocs.yml
                    # Configuração do mkdocs.
requirements.txt
                   # Bibliotecas python necessárias.
                    # Arquivo make com comandos úteis
.gitignore
                    # Configuração git.
README.md
                    # Descrição do projeto
docs/
                                # Documentação da estrutura do projeto.
    01_cep_cleaning.md
                                # Documentação do script de limpeza de CEP.
    02 geo search.md
                                # Documentação do script de busca das geolocalizações.
   03_distance_calculate.md
                                # Documentação do script para cálculo das distâncias.
    funs_cep.md
                                # Documentação do script com funções úteis.
    faixacep.txt
                   # Contém as faixas de CEPs e sua respectiva localidade. Construído a partir dos sites do correios.
    key_google.txt # Chave de acesso para a API do google. Solicitar ao autor.
data cep/ # Diretório com os **data frames** necessários para a execução do projeto. Solicitar ao autor.
sample/
   01 cep cleaning.py
                                # Limpeza de CEP.
                                # Busca de geolocalizações.
   03_distance_calculate.py
                               # Cálculo de distâncias
                                # Funções úteis.
    funs_cep.py
```

1.2 Para o projeto

Para poder atender aos requisitos necessários do projeto, escolha um dos tópicos a seguir que encaixe com as suas ferramentas:

Execute o arquivo Makefile com o GNU Make. Nele será criado um ambiente virtual e instalará todos os pacotes necessários do projeto. É necessário ter o virtualenv instalado.

Para executar via GNU Make:

```
make install
```

Caso não tenha o GNU Make. É aconselhável criar um ambiente virtual a partir da sua ferramente de escolha e execute o requirements.txt através do PIP.

Algumas ferramentas para a criação de ambientes virtuais:

```
virtualenv
venv
poetry
anaconda
```

Para executar via PIP:

```
pip install -r requirements.txt
```

1.3 Links úteis

Site de consulta de faixas de CEPs: https://buscacepinter.correios.com.br/app/faixa cep uf localidade/index.php

GNU Make: https://www.gnu.org/software/make/

Anaconda: https://www.anaconda.com

- 3/12 - ScientifiCloud

2. Script - Limpeza e formatação de CEPs

O script tem como objetivo limpar e formatar os CEPs para um padrão. O padrão adotado é para que o script de busca de geolocalizações possa ser o mais preciso possível.

A seguir será descrito as seções presentes no script:

2.1 Importando dados

Importa os dados com as informações de CEP a serem formatados.

2.2 Aplicando tratamento nos dados de CEP

Aplica a classe: FaixaCep, e as funções: limpar_cep(), formatar_cep(). Para mais informações, clique aqui.

2.2.1 Dicionário de faixas de CEP

Transforma o arquivo faixacep.txt em um dicionário para ser consultado.

```
dict_cep = eval(open("../config/faixacep.txt", "r").read())
faixacep = FaixaCep(dict_cep)
```

2.2.2 Tratamento

Limpa as strings de CEP, removendo qualquer caracter que não seja número e sinaliza como None os CEPs inválidos.

```
df["cepaluno"] = df["cepaluno"].astype("str")
df["cepaluno_limpo"] = df["cepaluno"].apply(limpar_cep)
df["cepunidade_limpo"] = df["cepunidade"].apply(limpar_cep)
```

O arquivo é filtrado para somente aqueles que contém CEP.

O processo é repetido para os alunos e para a unidade de estudo.

2.2.3 Formatação

Constrói a formatação a ser utilizada para a busca de geolocalização tanto para os alunos quanto para as unidades de estudo.

2.3 Exportando CEPs limpos

Exportando dados limpos para data cep/

2.3.1 Para mais informações

Para mais informações, consulte o texto Funções úteis.

- 4/12 - ScientifiCloud

3. Script - Busca de geolocalização

O script tem como objetivo buscar a geolocalização dos dados de Estado - CEP.

A seguir está descrito as seções presentes no script.

3.1 Importando dados limpos

Import os dados limpos obtidos do script de limpeza de CEPs.

3.2 Criando array apenas com CEPs distintos

Cria uma lista de CEPs únicos tanto de alunos quanto das unidades.

3.3 Buscando geolocalização

Primeiramente faz uma verificação para encontrar a base geocode.csv existente de geolocalizações para CEPs conhecidos.

```
if os.path.isfile(PATH):
    df_verify = pd.read_csv(PATH)
else:
    df_verify = pd.DataFrame(columns=column_name)
```

Faz um filtro para executar a API para os CEPs que não constam na base de dados. Caso não haja esse arquivo, é criado uma lista vazia e é feito a busca para todos os CEPs.

```
cep_distinct = [cep for cep in cep_distinct if cep not in df_verify["cep"].values]
```

Feito a lista de CEPs definitiva, uma lógica de lotes é executada. A API é executada em uma sequência de 100 CEPs. E após finalizado a sequência é salvo no arquivo geocode.csv as informações coletadas. Esse passo é executado até percorrer a lista de todos os CEPs.

```
BUFFER_SIZE = 100
for i in range(0, n_distinct, BUFFER_SIZE):
    start_at = i
    end_at = min(start_at + BUFFER_SIZE, n_distinct)
    temp = []
    print(f"Buscando CEPs {start_at} - {end_at}")
    for cep in cep_distinct[start_at:end_at]:
        geocode = geo.locate_geocode(cep)
        temp.append([cep, geocode])
    df_geocode = pd.DataFrame(temp, columns=column_name)

if not os.path.isfile(PATH):
    df_geocode.to_csv(PATH, index=False)
else:
    df_geocode.to_csv(PATH, mode="a", header=False, index=False)
```

3.3.1 Para mais informações

Para mais informações, consulte o texto Funções úteis.

4. Script - Calcular distância das geolocalizações

O script tem como objetivo executar o cálculo de distância entre duas geolocalizações.

A seguir será descrito as seções presentes no script:

4.1 Importando dados

Importa os dados contendo as informações das geolocalizações e os dados contendo os CEPs dos alunos e unidades.

4.2 Dando merge com os CEPs conhecidos

A partir dos dados contendo as geolocalizações e os CEPs únicos, é dado um merge com os CEPs de alunos.

```
df_aluno_geocode = df_notnull[["DR_CD_MATRICULA", "cepaluno_f"]].merge(
    df_unique_geocode, left_on="cepaluno_f", right_on="cep"
)
```

O mesmo é feito para os CEP de unidade:

```
df_unidade_geocode = df_notnull[["DR_CD_MATRICULA", "cepunidade_f"]].merge(
    df_unique_geocode, left_on="cepunidade_f", right_on="cep"
)
```

Contendo um banco de dados com as geolocalizações dos alunos e outro com as das unidades, é feito um merge com os dois bancos.

```
df_merged = df_aluno_geocode.merge(
    df_unidade_geocode,
    left_on="DR_CD_MATRICULA",
        right_on="DR_CD_MATRICULA",
    how="inner",
)

df_merged = df_merged.merge(
    df_notnull[["DR", "DR_CD_MATRICULA"]],
    left_on="DR_CD_MATRICULA",
        right_on="DR_CD_MATRICULA",
)
```

4.3 Cálculo de distância

Seção destinada para calcular a distância entre as geolocalizações

4.3.1 Filtrando para geolocalizações existentes

Remove as linhas que não contém informação da geolocalização do aluno ou da unidade.

4.3.2 Distância

É feito o cálculo da distância para a métrica de Haversine e para o GeoDesic. Feito os cálculos, um novo data frame é salvo com as informações:

- df distance full.csv Possui todas as colunas existentes;
- df_distance.csv Possui as colunas de DR_CD_MATRICULA, distance_hav e distance_dg.

4.3.3 Para mais informações

Para mais informações, consulte o texto Funções úteis.

5. Script - Funções úteis

5.1 FaixaCep

Bases: dict

Consulta de faixas de CEP de acordo com a sua localidade.

Parameters:

Name	Туре	Description	Default
dict	dict	Objeto contendo o dicionário com as faixas de CEP e sua localidade	required
		correspondente. O dict deve ser no formato: {range(start,end):	
		"localidade"}	

```
urce code in sample/funs_cep.py V
```

5.1.1 <u>__getitem__(item)</u>

Busca localidade a partir de uma consulta em faixas de valores.

Parameters:

Name	Туре	Description	Default
item	int	CEP a ser consultado.	required

Returns:

Name	Туре	Description
str	str	Localidade correspondete ao CEP inserido.

- 7/12 - ScientifiCloud

urce code in sample/funs_cep.py ~

```
def __getitem__(self, item: int) -> str:
    """Busca localidade a partir de uma consulta em faixas de valores.

Args:
    item (int): CEP a ser consultado.

Returns:
    str: Localidade correspondete ao CEP inserido.

"""

in the inistance(item, range):
    for key in self:
        if item in key:
        return self[key]
    return super().__getitem__(item)
```

5.2 GeoCode

Acessa a API do GoogleMaps para recuperar informações de latitude e longitude.

urce code in sample/funs_cep.py >

```
class GeoCode:
    """Acessa a API do GoogleMaps para recuperar informações de latitude e longitude."""
 91
92
93
94
95
             def __init__(self, path="../config/key_google.txt"):
    """Utiliza a client key no arquivo key_google.txt para estabelecer conexão com a API.
 96
97
                  ...g..
path (str): Diretório e nome do arquivo que está a key para a API google.
98
99
                  f = open(path, "r")
f = f.read()
100
101
                  self.key = eval(f)
self.gmaps_key = googlemaps.Client(key=self.key)
102
103
104
105
             def locate_geocode(self, cep: str) -> tuple:
                   """Consulta a latitude e longitude do CEP informado.
106
                  Args:
cep (str): CEP formatado, por exemplo: Distrito Federal - XXXXXX-XXXX. A função formatar_cep() formata o CEP da forma correta a ser utilizado.
107
108
                  Returns:
110
                  tuple: Latitude e Longitude a partir do CEP.
111
112
                  print(f"Buscando geolocalização para {cep}")
g = self.gmaps_key.geocode(cep)
113
114
115
116
                     self.lat = g[0]["geometry"]["location"]["lat"]
self.lon = g[0]["geometry"]["location"]["lng"]
return self.lat, self.lon
                  except:
                       return None
```

5.2.1 __init__(path='../config/key_google.txt')

Utiliza a client key no arquivo key_google.txt para estabelecer conexão com a API.

Parameters:

Name	Туре	Description	Default
path	str	Diretório e nome do arquivo que está a key para a API google.	<pre>'/config/ key_google.txt'</pre>

- 8/12 - ScientifiCloud

g3 def __init__(self, path="../config/key_google.txt"): 94 """Utliza a client key no arquivo key_google.txt para estabelecer conexão com a API. 95 96 Args: 97 path (str): Diretório e nome do arquivo que está a key para a API google. 98 """ 99 f = open(path, "r") 100 f = f.read() 101 self.key = eval(f) 102 self.gmaps_key = googlemaps.Client(key=self.key)

5.2.2 locate_geocode(cep)

Consulta a latitude e longitude do CEP informado.

Parameters:

Name	Туре	Description	Default
сер	str	CEP formatado, por exemplo: Distrito Federal - XXXXX-XXX. A função formatar_cep() formata o CEP da forma correta a ser utilizado.	required

Returns:

Name	Туре	Description
tuple	tuple	Latitude e Longitude a partir do CEP.

```
def locate_geocode(self, cep: str) -> tuple:
    """Consulta a latitude e longitude do CEP informado.

Args:
    cep (str): CEP formatado, por exemplo: Distrito Federal - XXXXX-XXXX. A função formatar_cep() formata o CEP da forma correta a ser utilizado.

Returns:
    tuple: Latitude e Longitude a partir do CEP.
    """

print(f"Buscando geolocalização para {cep}")
    g = self.gmaps_key.geocode(cep)
    try:
    self.lat = g[0]["geometry"]["location"]["lat"]
    self.lon = g[0]["geometry"]["location"]["lng"]
    return self.lat, self.lon
except:
    return None
```

5.3 calcular_distancia(origin, destination, metric='hav')

Cálculo de distância entre dois pontos a partir das suas geolocalizações.

Parameters:

Name	Туре	Description	Default
origin	tuple	Latitude e longitude da localização de origem.	required
destination	tuple	Latitude e longitude da localização de destino.	required
metric	str, opcional	Métrica a ser utilizada para o cálculo da distância. As métricas disponíveis são: hav - Haversine e gd - GeoDesic. Padrão "hav".	'hav'

- 9/12 - ScientifiCloud

Raises:

Туре	Description
KeyError	Métrica de cálculo errado. As opções válidas são: hav - Haversine ou gd - GeoDesic

Returns:

Name	Туре	Description
float	float	Distância em KM.

```
urce code in sample/funs_cep.py ~
        def calcular_distancia(
    origin: tuple, destination: tuple, metric: metrics = "hav"
126
127
128
129
130
131
       ) -> float:
"""Cálculo de distância entre dois pontos a partir das suas geolocalizações.
                  origin (tuple): Latitude e longitude da localização de origem.
destination (tuple): Latitude e longitude da localização de destino.
132
133
134
135
                  metric (str, opcional): Métrica a ser utilizada para o cálculo da distância. As métricas disponíveis são: hav - Haversine e gd - GeoDesic. Padrão
136
137
138
139
                   KeyError: Métrica de cálculo errado. As opções válidas são: hav - Haversine ou gd - GeoDesic
140
141
142
143
            Returns:
float: Distância em KM.
             if metric == "hav":
144
145
146
147
148
             distance = haversine(origin, destination)
elif metric == "gd":
    distance = GD(origin, destination).km
else:
```

5.4 formatar cep(cep, faixacep)

raise KeyError("Metric incorrect!")

Formatação do CEP com a localidade.

return distance

Concatena a informação da localidade (UF) e o CEP.

Parameters:

Name	Туре	Description	Default
сер	str or int	CEP.	required
faixacep	FaixaCep	Objeto FaixaCep contendo as faixas de CEP e sua localidade correspondente.	required

Returns:

Name	Туре	Description
str	str	CEP formatado, exemplo: Distrito Federal - XXXXXXXX.

- 10/12 - ScientifiCloud

urce code in sample/funs_cep.py >

```
def formatar_cep(cep, faixacep: FaixaCep) -> str:
    """Formatação do CEP com a localidade.

73
74    Concatena a informação da localidade (UF) e o CEP.
75
76    Args:
77     cep (str or int): CEP.
78     faixacep (FaixaCep): Objeto FaixaCep contendo as faixas de CEP e sua localidade correspondente.
79
80    Returns:
81     str: CEP formatado, exemplo: Distrito Federal - XXXXXXXXX.
82    """
83
84    estado = faixacep[cep]
85    cep_formatado = f"(estado) - {str(cep)}"
86    return cep_formatado
```

5.5 limpar_cep(cep)

Limpa strings de CEP.

Retira caracteres que não sejam números. Retira CEPs com tamanhos incorretos, com 1 dígito em toda a string e retira aqueles que iniciam em uma faixa inexistente.

As faixas consultadas estão no arquivo faixas_cep.txt. As remoções existentes são aquelas que iniciam com:

remove_tuple = ("789", "76", "75")

Parameters:

Name	Туре	Description	Default
сер	str	CEP	required

Returns:

Name	Туре	Description
int	int	CEP em formato de int para que possa ser consultado no dicionário de faixas de CEP. Consultar documentação de FaixaCep().

- 11/12 - ScientifiCloud

urce code in sample/funs_cep.py ~

```
def limpar_cep(cep: str) -> int:

"""Limpa strings de CEP.

Retira caracteres que não sejam números. Retira CEPs com tamanhos incorretos, com l dígito em toda a string e retira aqueles que iniciam em uma faixa inexistente.

As faixas consultadas estão no arquivo faixas_cep.txt. As remoções existentes são aquelas que iniciam com:

remove_tuple = ("789", "76", "75")

Args:
cep (str): CEP

Returns:
int: CEP em formato de int para que possa ser consultado no dicionário de faixas de CEP. Consultar documentação de FaixaCep().

"""

remove_tuple = ("789", "76", "75")
cep = re.sub("[0-9]", "", cep)
if len(cep) < 7 or len(cep) > 8:
cep = None
elif cep.count(cep[0]) == len(cep):
cep = None
elif cep.startswith(rmove_tuple):
cep = lone
elif cep.startswith(rmove_tuple):
cep = lone
elif cep.startswith("76"):
cep = lone
elif cep.startswith("76"):
cep = lone
if cep: startswith("76"):
cep = lone
return int(cep)
return cep
```