

Contents

1. Sử dụng Open Street Map - request	2
2. Sử dụng Selenium	4
3. Ghép file	8

1. Sử dụng Open Street Map - request

```
request.py > ...
1  import requests
2  import time
3  import pandas as pd
4
5  def get_places_of_interest(latitude, longitude):
6      url = f"https://nominatim.openstreetmap.org/reverse"
7      params = {
8          "format": "json",
9          "lat": latitude,
10         "lon": longitude,
11         "addressdetails": 1,
12     }
13
14     response = requests.get(url, params=params)
15     data = response.json()
16     places_of_interest = []
17     for key in ['building', 'hospital', 'factory', 'apartments', 'office', 'residential', 'aeroway']:
18         if key in data['address']:
19             places_of_interest.append(data['display_name'])
20             return places_of_interest
21         else:
22             continue
23
24
25
26
27 df = pd.read_csv("4g.csv")
28 address1 = []
29 latitude = df["LATITUDE"]
30 longitude = df["LONGITUDE"]
31 # Nhập vào tọa độ latitude và longitude
32 for i in range(0, len(latitude)):
33     print(i)
34     places_of_interest = get_places_of_interest(latitude[i], longitude[i])
35     if places_of_interest:
36         address1.append(places_of_interest[0])
37     else:
38         address1.append(places_of_interest)
39
40 print("_____")
41 print("Completed")
42 df["Address1"] = address1
43 # Tên tệp Excel là 'output.xlsx', index=False để không bao gồm cột index
44 df.to_excel('4g_output1_use_libRequest.xlsx', index=False)
45
```

ở dòng 17 để chúng ta có thể chọn được các vị trí mong muốn

chúng ta có thể tham khảo ở link này : <https://taginfo.openstreetmap.org/>

đầu vào là file 4g.csv

đầu ra là file 4g_output1_use_libRequest . các tọa độ mà ko lấy được thông tin về sẽ để Null

Sau khi lọc ta sẽ được 2 file :

4g_output1_use_request_win

4g_output1_use_request_fail

```
filter_request.py > ...
1  import pandas as pd
2  import re
3  # đọc dữ liệu từ file excel
4  df = pd.read_excel("4g_output1_use_libRequest.xlsx")
5  print(df.info())
6  df1 = df[df['Address1'].isna()]
7  print(df1.info())
8
9  df2 = df[df['Address1'].notna()]
10 print(df2.info())
11
12
13
14 df1.to_excel('4g_output1_use_request_fail.xlsx', index=False)
15 df2.to_excel('4g_output1_use_request_win.xlsx', index=False)
```

4g_output1_use_request_fail sẽ đi qua tool Selenium

2. Sử dụng Selenium

```
1  from selenium import webdriver
2  from selenium.webdriver.common.by import By
3  from selenium.webdriver.common.keys import Keys
4  from bs4 import BeautifulSoup
5  import time
6  from selenium.webdriver.chrome.options import Options
7  import pandas as pd
8  import pickle
9
10
11  try:
12      df = pd.read_excel("4g_output1_use_request_fail.xlsx")
13      address2 = []
14      latitude = df["LATITUDE"]
15      longitude = df["LONGITUDE"]
16      chrome_options = Options()
17      # chạy nền ko hiển thị lên
18      chrome_options.add_argument("--headless")
19      driver = webdriver.Chrome(options= chrome_options)
20      time.sleep(0.1)
21      driver.get("https://www.google.com/maps")
22      driver.implicitly_wait(20)
23      for i in range (950,len(latitude)):
24          print(i+1)
25          if(i%50 == 49):
26              time.sleep(1)
27              driver = webdriver.Chrome(options= chrome_options)
28              driver.get("https://www.google.com/maps")
29              driver.implicitly_wait(30)
30
31      # Tìm ô tìm kiếm bằng ID
32      search_box = driver.find_element(By.ID, "searchboxinput")
33      # Nhập tọa độ vào ô tìm kiếm
34      search_box.send_keys(f"{latitude[i]}, {longitude[i]}")
35      # Cách 1
36      # Tìm nút tìm kiếm bằng ID và bấm vào nó
37      #search_button = driver.find_element(By.ID, "searchbox-searchbutton")
38      #search_button.click()
39      # Cách 2 gửi yêu cầu truy cập trực tiếp bằng phím ENTER
40      driver.implicitly_wait(30)
41      search_box.send_keys(Keys.ENTER)
42
43      # Chờ cho kết quả xuất hiện (thời gian chờ có thể cần điều chỉnh)
44      driver.implicitly_wait(30)
45      time.sleep(0.1)
46      # Thu thập thông tin về địa chỉ từ trang tìm kiếm
47      place_address = driver.find_element(By.CSS_SELECTOR, "#QA0Szd > div > div > div.w6VYqd
```

Tìm ô searchbox để fill tọa độ vào

The image shows a Google Maps interface with the search bar highlighted. The search bar contains the text "Tìm kiếm trên Google Maps" and "ut.xiQnY". The search bar is a text input field. The DevTools console shows the following HTML structure for the search bar:

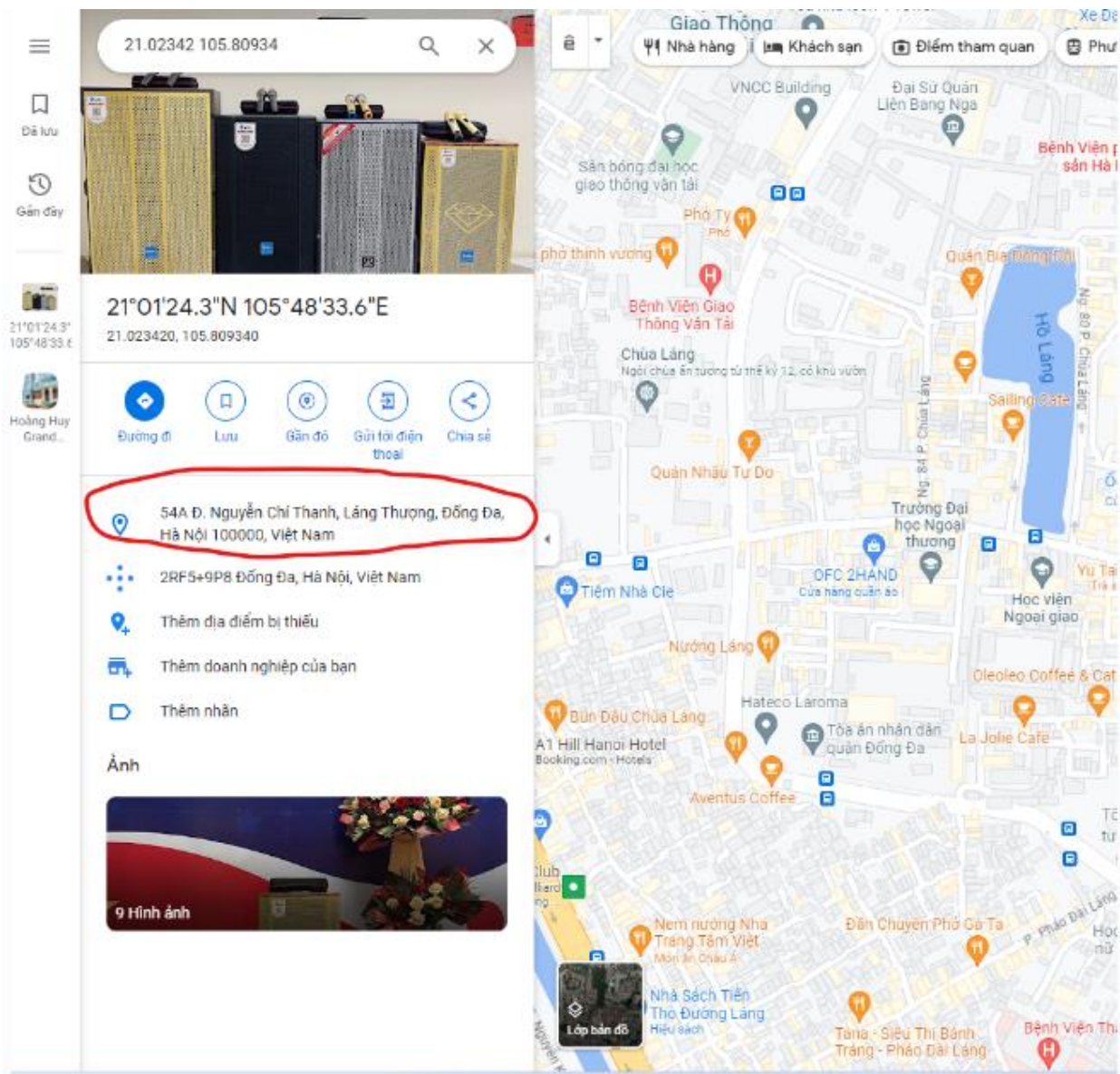
```

<div class="input searchboxinput.searchboxinp ut.xiQnY" data-bbox="114 86 884 649">
  <input type="text" value="Tìm kiếm trên Google Maps" />
</div>

```

The search bar is a text input field. The DevTools console shows the following HTML structure for the search bar:

Sau đó thông tin trả về như sau



Để lấy được ở ô màu đỏ ta sử dụng css_selector tìm đúng vị trí cần lấy
Sau đó lại nhập tọa độ mới


```

46 # Thu thập thông tin về địa chỉ từ trang tìm kiếm
47     place_address = driver.find_element(By.CSS_SELECTOR, "#QA0Szd > div > div > div.w6VVqd
48     driver.implicitly_wait(30)
49     clear_box = driver.find_element(By.CSS_SELECTOR, "#searchbox > div.lSDxNd > button")
50     clear_box.click()
51     driver.implicitly_wait(30)
52 # In thông tin địa chỉ
53     address2.append(place_address)
54 # Đóng trình duyệt
55     driver.quit()
56
57 except Exception as e:
58     # Xử lý lỗi (nếu cần)
59     print("Xảy ra lỗi:", str(e))
60
61 finally:
62     # Lưu danh sách thông tin thu thập được vào file pickle
63     with open('address22.pkl', 'wb') as f:
64         pickle.dump(address2, f)
65
66
67 print("_____")
68 print("Completed")
69 df["Address2"] = address2
70 # Tên tệp Excel là 'output.xlsx', index=False để không bao gồm cột index
71 df.to_excel('4g_output2_use_selenium.xlsx', index=False)

```

ở đoạn sau chúng ta sử dụng finally: vì chạy selenium hay bị đứt đường truyền . Đứt ở đâu thì lưu những thông tin đã thu thập được để khi chạy lại biết chạy từ chỗ nào

Đầu vào là file 4g_output1_use_request_fail

Đầu ra là file đã fill đầy đủ thông tin về địa chỉ (những có cả tòa nhà cả tên đường ...)

Ta sẽ sử dụng filter để lọc các giá trị mong muốn

```

filter_selenium.py > ...
1 import pandas as pd
2 import re
3 # đọc dữ liệu từ file excel
4 df = pd.read_excel("4g_output2_use_selenium.xlsx")
5 #Danh sách các từ bạn muốn tìm kiếm trong địa chỉ
6 desired_words = ['khu đô thị', 'tower', 'tòa', 'chung cư', 'building', 'vin', 'plaza', 'land', 'home', 'apartment', 'chung cư', 'căn hộ', 'center', 'nhà', 'KĐT',
7                 'big c', 'văn phòng', 'bệnh viện', 'hotel', 'grand', 'khách sạn', 'trường', 'mall', 'bank', 'sun', 'garden', 'park', 'trung tâm', 'sky',
8                 'pearl', 'công ty', 'ct', 'hh', 'win', 'house', 'town', 'holding', 'Riverside', 'department', 'mart', 'ngân hàng', 'cao ốc', 'Appartments', 'Khu thương mại',
9                 'Tòa Nhà', 'Apartments', 'home', 'toa nha', 'celadon', 'block', 'gele', 'office', 'green', 'Tòa nhà' ]
10
11 # Tạo biểu thức chính quy với các từ bạn muốn tìm kiếm, kết hợp chúng bằng toán tử |
12 regex_pattern = '|'.join(map(re.escape, desired_words))
13 # Lọc các hàng trong đó cột 'Address1' chứa ít nhất một từ trong danh sách
14 filtered_df_win = df[df['Address1'].str.contains(regex_pattern, case=False, na=False, regex=True)]
15 filtered_df_fail = df[~df['Address1'].str.contains(regex_pattern, case=False, na=False, regex=True)]
16 print(len(filtered_df_win))
17 print(len(df))
18 print(len(filtered_df_win)/len(df))
19 print(len(filtered_df_fail))
20
21 filtered_df_win.to_excel('4g_output2_use_selenium_win.xlsx', index=False)
22 filtered_df_fail.to_excel('4g_output2_use_selenium_fail.xlsx', index=False)

```

Sau khi Filter xong ta sẽ thu được 2 file .

1 file thỏa mãn yêu cầu lọc các các từ như chúng ta đã đặt ra : 4g_output2_use_seleniumm_win

1 file ko thỏa mãn : 4g_output2_use_seleniumm_fail

3.Ghép file

Ghép 1 file thỏa mãn ở phần sử dụng Open street map (4g_output1_use_request_win)và file thỏa mãn của phần selenium (4g_output2_use_seleniumm_win)

Ta được 1 file các tọa độ trả về đúng tên tòa nhà như yêu cầu (4g_win).ta sẽ biết được số lượng hàng trong file này là bao nhiêu để tí còn điền cột “status”

Ghép tiếp file vừa nhận được ở trên (4g_win). với file (4g_output2_use_seleniumm_fail) ta được file hoàn chỉnh (4g_complete). đủ số lượng bản ghi như file ban đầu (file 4g ban đầu được nhận là 3838 dòng, thì file này cũng đủ số lượng 3838 dòng như thế)

Tiếp theo thêm cột “status”

```
status.py > ...
1  import pandas as pd
2
3  df = pd.read_excel("4g_win.xlsx")
4  length_4g_win = len(df)
5
6
7  df1 = pd.read_excel("4g_complete.xlsx")
8  df['status'] = 0 # Khởi tạo tất cả giá trị là 0
9  df['status'].iloc[:length_4g_win] = 1 # Gán giá trị 1 cho length hàng đầu tiên
10 print(df)
11 df.to_excel('4g_completeddd.xlsx', index=False)
12 |
```

Đầu vào file 4g_complete

Đầu ra là file 4g_completeddd đã được đánh trạng thái 0 1 (1 tức là ra được tên tòa nhà , 0 tức là chỉ ra được địa chỉ)

ở dòng thứ 4 chính là số dòng của file 4g_win