

```

1  #!/usr/bin/env python
2  # coding: utf-8
3  # Import required Python libraries
4  import requests
5  from bs4 import BeautifulSoup, re, Comment
6  import pandas as pd
7  import xlswriter
8
9  # Step 1: Scraping the primary "MEPS data file website",
10 # finding the data file names that are within the "option"
11 # comment tags, and saving them in a csv file
12
13 def extractOptions(inputData):
14     sub1 = str(re.escape('<option value="All">All data files</option>'))
15     sub2 = str(re.escape('</select>'))
16     result = re.findall(sub1+"(.*)" + sub2, inputData, flags=re.S)
17     if len(result) > 0:
18         return result[0]
19
20 def extractData(inputData):
21     sub1 = str(re.escape('>'))
22     sub2 = str(re.escape('</option>'))
23     result = re.findall(sub1+"(.*)" + sub2, inputData, flags=re.S)
24     if len(result) > 0:
25         return result[0]
26     return ''
27
28 def main(base_url):
29     response = requests.get(base_url)
30     soup = BeautifulSoup(response.text, "html.parser")
31     comments = soup.find_all(string=lambda text: isinstance(text, Comment))
32
33     for c in comments:
34         if '<select id="pufnumber" size=1 name="cboPufNumber">' in c:
35             options = extractOptions(c)
36             ops = options.splitlines() #split text into lines
37             fp = open(r'C:/Data/MEPS_fn.csv', 'w')
38             for op in ops:
39                 data = extractData(op)
40                 if data != '': #check if the data found
41                     fp.write(data + '\n')
42             fp.close()
43
44             with open(r'C:/Data/MEPS_fn.csv', 'r') as buff:
45                 for i, line in enumerate(buff, 1):
46                     pass
47                 print(f"({i})", 'file names listed in the MEPS website')
48
49 main('https://meps.ahrq.gov/data_stats/download_data_files.jsp')
50
51 # Step 2: Creating a Pandas DataFrame from the csv file 12/31/2022
52
53 colname = ['file_name']
54 df1 = pd.read_csv(r'C:/Data/MEPS_fn.csv', sep='\t', names = colname)
55
56 df1.drop(df1[df1['file_name'].str.contains('replaced|CD-ROM|NHC|NHEA|NHIS Link|HC-IC Linked|
1996 Parent IDs')].index, inplace=True)
57
58 df1["file_id"] = df1["file_name"].str.extract(r"([A-Z])[A-Z]+-(\d+[A-
Z]*)").sum(axis=1).str.lower()
59 df1['file_id'] = df1['file_id'].str.replace('h0', 'h').str.replace('h36', 'h036')
60 .str.replace('h36brr', 'h036brr')
61
62 df1["url1"] =
63 "https://meps.ahrq.gov/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-" +
64 df1["file_name"].str.extract(r"(\d+[A-Z]*)").sum(axis=1).astype(str)
65
66 df1.reset_index(drop = True, inplace = True)
67
68 print("{:,}".format(len(df1)), 'MEPS public-use file names')
69

```

```

67 # Step 3: Scraping all MEPS data file-specific websites
68 # and saving format-specific file names from each of
69 # those sites in a DataFrame 12/31/2022
70
71 url2_str_list = []
72 for item in df1.index:
73     url1_str = df1['url1'][item]
74     response = requests.get(url1_str)
75     soup = BeautifulSoup(response.text, "html.parser")
76
77     for link in soup.find_all('a'):
78         if link.text.endswith('.ZIP'):
79             url2_str = 'https://meps.ahrq.gov' + link.get('href').strip('.')
80             url2_str_list.append(url2_str)
81
82 df2 = pd.DataFrame(url2_str_list, columns=['url2'])
83 df2['file_id'] = df2['url2'].str.extract(r"([h]\d+[abcdefghijklmnopqrstuvwxyz]*?!\\d)").sum(axis=1)
84 df2['file_id'] = df2['file_id'].str.replace('da', '')
85
86 df1 = df1.drop('url1', axis=1)
87 merged_df = pd.merge(df1, df2, on='file_id', validate="one_to_many")
88 print("{:,}".format(len(merged_df)), 'URLs that are specific to data file formats')
89 with pd.ExcelWriter('merged_df.xlsx') as writer:
90     merged_df.to_excel(writer, sheet_name='data_urls', index=False)
91     writer.sheets['data_urls'].set_column(45, 3, 45)
92

```