# One Click to Create a List of All MEPS-HC ZIP Files' Download URLs: Web Scraping with SAS® and Python

Pradip K. Muhuri[1], Agency for Healthcare Research and Quality
Charles Z.X. Han[1], Ryman Healthcare Ltd.
John Vickery[1], North Carolina State University

## ABSTRACT

This paper presents a case study of web scrapping the Medical Expenditure Panel Survey-Household Component (MEPS-HC) sites with SAS® and Python. Each program dynamically constructs an up-to-date list of over one thousand clickable MEPS-HC ZIP files' download links (Uniform Resource Locators or URLs) and several identifiers and then saves the list into an Excel file. Once you automate the creation of a list of all those URLs in an Excel file using one of these programs (SAS or Python), you can directly initiate ZIP file downloads by simply clicking the relevant URLs from that file outside of the Internet, an easy alternative to navigating multiple pertinent websites for the same task. These programs are reusable to obtain an updated list of MEPS-HC file download links as more data sets are added to the website annually from the ongoing panel survey.

## INTRODUCTION

The Agency for Healthcare Research and Quality (AHRQ) has made various data files available from the ongoing Medical Expenditure Panel Survey-Household Component (MEPS-HC[2]) to the public on its website. On this site (https://meps.ahrq.gov/data_stats/download_data_files.jsp), you can manually select an individual data year or all available years from a dropdown list and check the box that represents the file type(s)[3] for your search and then navigate through the site for the desired public-use file (PUF). The PUF site includes links for the data file documentation, survey questionnaires, codebooks, SAS®/Stata/R programs to read the raw data into the individual file format, and ZIP files' download URLs.

Over one thousand ZIP files' download URLs for ASCII, Excel, SAS transport, SAS V9, and Stata files exist across all MEPS-HC PUF websites. However, a consolidated list of all those URLs is unavailable at a single location. Therefore, finding an HTML link to trigger a ZIP file download can be tedious and time-consuming when desiring many such files because the process requires manual navigation across multiple linked websites.[4] The inefficiency increases with the number of file downloads. The goal of web

---

[1] Pradip K. Muhuri conceived the idea of the paper, drafted it, and wrote the earlier version of the SAS and Python programs in it. Charles Z.X. Han subsequently revised the Python program and contributed to the code description. Later, John Vickery optimized the SAS and Python programs to their current form. Finally, Pradip K. Muhuri worked on another round of manuscript revisions, to which Charles Z.X and John Vickery contributed.

[2] MEPS is a set of large-scale on-going surveys of families and individuals, their medical providers, and employers across the United States. It is the complete data source on the cost and use of health care and health insurance coverage for the U.S. civilian noninstitutionalized population [1]. See https://www.meps.ahrq.gov/mepsweb/about_meps/survey_back.jsp for a complete description of MEPS, including its Household Component (HC).

[3] Examples of file types include Household Full Year File, Household Event File (e.g., office-based medical provider visits file, outpatient visits file, emergency room visits file, hospital in-patient stays file, home health visits file, dental visits file, and prescribed medicines file), and Pooled Linkage File.

[4] For the data year 2023, the plans spread over ten months from March through December (https://meps.ahrq.gov/mepsweb/about_meps/releaseschedule.jsp).

scraping the MEPS-HC sites with SAS® and Python here is to dynamically construct an up-to-date list of zip files' download URLs and several identifiers and save the list in individual Excel files for easy use.

## SAS® PROGRAM

We present the SAS program[5] (shown in Appendix I) in six steps below.

**Step 1: Fetch the content of the main MEPS-HC web page**

**Lines 3-7**: The code performs the HTTP requests to the specified URL. It writes the response body into the fileref SOURCE[6].

**Step 2: Extract data year values by parsing HTML content**

**Line 11:** The DATA statement creates an out**put** SAS data set YEAR_VALUES.

**Line 13**: The INFILE statement specifies the fileref SOURCE and other relevant options (e.g., LENGTH=, LRECL=, and END=).

SAS PRX (Perl Regular Expression) functions and call routines are used to extract the data year value from within the HTML option tags (examples below). As of this writing, there are 26 MEPS PUF data years, ranging from 1996 to 2021.

```
<option value="2021">2021</option>
<option value="2020">2020</option>
…
 <option value="1996">1996</option>
```

**Line 14**: The PRXPARSE function compiles Perl regular expressions in the DATA step. It defines and retains the pattern in the numerical identifier variable RE.

```
re = prxparse('/<option value="\d{4}">(\d{4})<\/option>/');
```

The <option value="\d{4}"> is just matching the HTML option tag code, but the \d{4} within the parentheses is the 4-digit value that is being captured.

**Line 16**: The DO WHILE with the test condition NOT EOF in parentheses iterates while the condition is true, executing the subsequent statements in the same code block.

**Line 17**: The INPUT statement brings the data record into the input buffer, creating a SAS variable HTML_LINE by specifying the INFORMAT $VARYING32767 and a required numeric variable RECLEN to it.

---

[5] We used SAS® 9.04.01M6 X64_10PRO in the Windowing environment.
[6] Immediately above PROC HTTP, the FILENAME statement associates the file reference (FILEREF) named SOURCE with the external file (directory and filename).

**Lines 19-23**: The code performs a logical test using the PRXMATCH function[7] with the IF statement.

```
if prxmatch(re, html_line) > 0 then do;
   call prxposn(re, 1, start, end);
   year = substr(html_line, start, end);
   output year_values;
 end;
```

This IF statement tests the condition of whether the starting position of the match using PRXMATCH is greater than 0. If there is a match, three things happen. First, the call to PRXPOSN uses the results (RE) from the PRXPARSE to return a capture buffer (i.e., 1 in this case)[8]. Then, it provides the results to the SUBSTR function to return the text (i.e., the value of the variable YEAR) from the string HTML_LINE. Finally, the OUTPUT statement writes the observation to the data set YEAR_VALUES.

The partial SAS Log:

```
NOTE: 3244 records were read from the infile SOURCE.
      The minimum record length was 0.
      The maximum record length was 452.
NOTE: The data set WORK.YEAR_VALUES has 26 observations and 5 variables.
```

Below is the PROC PRINT output from WORK.VALUES (5 observations).

| html_line | re | start | end | year |
|---|---|---|---|---|
| <option value="2021">2021</option> | 1 | 30 | 4 | 2021 |
| <option value="2020">2020</option> | 1 | 30 | 4 | 2020 |
| <option value="2019">2019</option> | 1 | 30 | 4 | 2019 |
| <option value="2018">2018</option> | 1 | 30 | 4 | 2018 |
| <option value="2017">2017</option> | 1 | 30 | 4 | 2017 |

**Lines 28-30**: The NODUPKEY option with PROC SORT statement overwrites the output data set YEAR_VALUES with no duplicate observations (i.e., only one observation for the _ALL_ variable value). Note the KEEP= data set option listing the variable YEAR after the data set name within the same procedure statement. This option causes the output data set to hold the variable YEAR only. .

The Partial SAS Log:

```
NOTE: There were 26 observations read from the data set WORK.YEAR_VALUES.
NOTE: 0 observations with duplicate key values were deleted.
NOTE: The data set WORK.YEAR_VALUES has 26 observations and 1 variables.
```

---

[7] The first argument of the PRXMATCH function is the variable RE (whose value is the regular expression) created by the PRXPARSE function, while the second argument is the name of the character string (HTML_LINE) to search for or match against. PRXMATCH returns the first position of a match or 0 if there is no matching.
[8] The number of parentheses pairs in the regular expression determines the value that identifies the capture buffer from which to obtain the start position and length.

**Step 3: Call PROC HTTP and extract public-use file numbers and year URLs via macro processing**

**Lines 33-88**: The macro GET_YEAR with the parameter YEAR wraps PROC HTTP and DATA step in the macro definition. The CALL EXECUTE routine in the DATA _NULL_ step generates the macro calls. With the execution of this macro, the macro variables get resolved. As examples, the resolved values of two macro variables from the first iteration are below.

- &YEAR resolves to 1996
- &YEAR_URL resolves to
  https://meps.ahrq.gov/data_stats/download_data_files_results.jsp?cboDataYear=1996&buttonYearandDataType=Search

Thus the macro GET_YEAR generates text for 26 PROC HTTP and 26 DATA steps.

Within the macro definition, the DATA step uses the PRXPARSE function plus CALL PRXPOSN and the SUBSTR function following a logical test using the PRXMATCH function with the IF statement to extract values for the PUF number (PUF_NUM), the MEPS file name (MEPS_FILE), and the data year (DATA_YEAR).

**Line 69:** Some search results include strings like <br> tags in the value of the variable MEPS_FILE (code and output below).

```
meps_file = substr(html_line, start, end);
```

The output: `2020 Full Year Population Characteristics File<br/>(HC-219 replaced by HC-224)`

**Line 70:** Stary HTML tags in the search results above require removal. So, the PRXCHANGE function removes undesired tags in the value of the variable MEPS_FILE.

```
meps_file = prxchange('s/<.+>//', -1, meps_file);
```

The output: `2020 Full Year Population Characteristics File(HC-219 replaced by HC-224)`

Of the various SAS data sets created, the PROC PRINT output from WORK.YEAR_LIST_2021 is below.

```
                                               data_
 PUF_num     meps_file                          year

 HC-228      2021 Full Year Population Characteristics File   2021
 HC-227      2021 Jobs File                      2021
```

**92-95:** Note the  SAS data set names ending in numeric year suffixes (e.g., YEAR_LIST_1996, YEAR_LIST_1997. YEAR_LIST_2021). The DATA step concatenates the SAS data sets using the name prefix (YEAR_LIST) followed by a colon(**:**) to select those ending in numeric suffixes (i.e., YEAR values)

in the SET statement. It then filters the data sets whose PUF_NUM begins with 'HC'[9]. The resultant output data set is WORK.YEAR_LIST.

**Lines 98-100**: The NODUPKEY option in the PROC SORT statement creates an output data set YEAR_VALUES with no duplicate observations (i.e., only one observation for the _ALL_ variable values).

The partial SAS Log[10]:

```
NOTE: There were 509 observations read from the data set WORK.YEAR_LIST.
NOTE: 79 observations with duplicate key values were deleted.
NOTE: The data set WORK.YEAR_LIST has 430 observations and 3 variables.
NOTE: PROCEDURE SORT used (Total process time):
```

**Step 4: Call PROC HTTP and extract public-use file (PUF) URLs to construct ZIP file URLs via macro processing**

**Lines 103-150**: The macro GET_PUF with one positional parameter (PUFNUM) wraps PROC HTTP and DATA step in the macro definition. The CALL EXECUTE routine in the DATA _NULL_ step generates the macro calls. With the execution of this macro, the macro variables get resolved. As examples, the resolved values of three macro variables from the first iteration are below.

- &PUFNUM resolves to HC-092
- &PUF_URL resolves to
  https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC092
- PUF_LIST_%sysfunc(translate(&PUFNUM, '_', '-')) resolves to PUF_LIST_HC_092

This macro GET_PUF generates text for hundreds of URL schemes to create HTTP requests and over one thousand DATA steps to construct  ZIP files' download URLS.

Within the macro definition, the DATA step uses the PRXPARSE function plus CALL PRXPOSN and the SUBSTR function following a logical test using the PRXMATCH function with the IF statement to extract values for the PUF number (PUF_NUM) and the file format (FILE_FORMAT). Note using the CATX  and SUBSTR functions to create the ZIP file link (ZIP_LINK).

---

[9] PUF_NUMs not beginning with 'HC' include 26 NHIS Link Files, 15 NHEA-Aligned MEPS Files, and eleven NHC_: and LINK_: files (the colon is a wild character) as of this writing.

[10] The duplicate values include two or more occuurences of the MEPS Longitudinal Data File from various panels, the 2000-2013 Employment Variables  File, the 2002-2009 Risk Adjustment Scores File, the 1996-2001 Risk Adjustment Scores File (HC-081 replaced by HC-092), the 2001 and 2002 MEPS HC Survey Data (CD-ROM), the Multum Lexicon Addendum Files to MEPS Prescribed Medicines Files 1996-2013, and the 1999 and 2000 MEPS HC Survey Data (CD-ROM).

The next is the PROC PRINT output from WORK.ZIP_LIST_HC_228.

```
PUF_num    file_format                                zip_link


HC-228    ASCII format            https://meps.ahrq.gov/data_files/pufs/h228/h228dat.zip
HC-228    SAS transport format    https://meps.ahrq.gov/data_files/pufs/h228/h228ssp.zip
HC-228    SAS V9 format           https://meps.ahrq.gov/data_files/pufs/h228/h228v9.zip
HC-228    Stata format            https://meps.ahrq.gov/data_files/pufs/h228/h228dta.zip
HC-228    XLSX format             https://meps.ahrq.gov/data_files/pufs/h228/h228xlsx.zip
```

**Lines 147-150**: The DATA step concatenates the SAS data sets[11] using the name prefix (PUF_LIST) followed by a colon(**:**) to select those ending PUF numbers as suffixes. The resultant output data set is WORK.PUF_LIST.

**Lines 158-160**:  The NODUPKEY option in the PROC SORT statement creates an output data set PUF_LIST with no duplicate observations (i.e., only one observation for the _ALL_ variable values).

The partial SAS Log:

```
NOTE: There were 1179 observations read from the data set WORK.PUF_LIST.
NOTE: 0 observations with duplicate key values were deleted.
NOTE: The data set WORK.PUF_LIST has 1179 observations and 3 variables.
```

**Step 5: One-to-many matching  - join up extracted pieces next to each other**

We combine extracted data elements horizontally. It is like joining knitted pieces.

**Lines 163-170**: The PROC SQL performs a one-to-many match with the tables YEAR_LIST and PUF_LIST based on the key column PUF_NUM and creates a table MEPS_ZIP_LINKS. Note the following.

- The table YEAR_LIST has a unique value for the column PUF_NUM.
- The table PUF_LIST includes multiple rows with the same value as the column PUF_NUM.
- The output table MEPS_ZIP_LINKS has a unique value for the column ZIP_LINK.

The SAS Log:

```
NOTE: Table WORK.MEPS_ZIP_LINKS created, with 1179 rows and 5 columns.
```

**Step 6:  Create listings with PROC REPORT and output them into an Excel spreadsheet**

**Lines 177-194:**  The REPORT procedure identifies the final merged data set, displaying the values of all specified variables. The CALL DEFINE associates the URL with ZIP_LINK. The ODS EXCEL statement

---

[11] Over 20 SAS data sets created here are empty because the corresponding ZIP files no longer existed on the MEPS website.  They were actually replaced with different PUF numbers. For example, HC-002 was replaced by HC-012  See https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=HC-002.

directs the PROC REPORT output to a file listed in the FILE= option. With the execution of the ODS EXCEL CLOSE; statement, the Excel file gets created (part of the output from the Excel spreadsheet shown in APPENDIX II).

The combined partial SAS Log:

```
NOTE: There were 1179 observations read from the data set WORK.MEPS_ZIP_LINKS
NOTE: Writing EXCEL file: c:\SESUG_2023\SAS_MEPS_zip_links_2023-06-07.xlsx
```

## PYTHON PROGRAM

We used Python (Version 3.9.12) in the Spyder IDE and JupyterLab Python Notebook (Version 6.4.5) as our second solution for web scraping. We present the Python program (as shown in APPENDIX III) in four steps below.

**Step 1: Import libraries**

**Lines 8-11:** The program imports necessary libraries, including Requests and BeautifulSoup[12] to scrape and parse MEPS-HC websites and dynamically create a list of over one thousand URLs that would trigger data downloads.

**Step 2: Get the list of numerical year options and saving as the list of tuples with year and URL**

**Lines 18-19**: This code uses the 'requests' library to send a GET request to a website with the URL "https://meps.ahrq.gov/data_stats/download_data_files.jsp". The response from the website is then passed to BeautifulSoup, another library, which parses the HTML text into an object called main_soup.

**Line 23**: The code defines a base URL as "https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_results.jsp?cboDataYear=", used to construct the URLs for downloading data files.

**Line 24-26**: The 'year_url_suffix' variable is set to "&buttonYearandDataType=Search", which will be appended to each year-specific URL. Next, the code extracts the available years for data download from 'main_soup'. It does so by finding the first 'select' element with the name attribute set to "cboDataYear", which is used to choose a specific year for data download. The '.select("option")' method is then called on this `select` element, which returns all the 'option' elements. The 'year_options' variable now contains a list of all the available years for data download.

The code below prints the list `year_options`.

```
print(year_options)
```

The output:

---

[12] BeautifulSoup, Scrapy, and Selenium are the three main Python-based tools commonly used for web scraping. BeautifulSoup() of the Python bs4 library enables one to extract data (specific elements) from a single webpage at a time effectively. Scrapy is a web scraping framework that can crawl various web pages, downloading, parsing, and storing data, whereas Selenium can automate navigating to sites and scrape the webpage content dynamically.

**Lines 28-31** The code creates a 'year_url_list' by iterating through each "option" element in 'year_options'. If the "value" attribute of the "option" element is a digit, a tuple is added to 'year_url_list'. The first element of the tuple is the value of the "value" attribute, representing the year for data download. The second element is a URL constructed by concatenating the base URL with the year value, passed through the 'quote()' method to ensure it is appropriately encoded for use in a URL. Finally, the 'year_url_suffix' variable is appended to the URL.

Overall, this code retrieves available data years associated with the public-use file name from the website, constructs URLs for downloading data files for each year, and stores them in a list of tuples called 'year_url_list'.

Below is the additional code (not shown in the Python program in Appendix III) to print the first two elements in the list `year_url_list`.

```
for item in year_url_list[:2]:    print(item)
```

The output:

```
('2021', 'https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_results.jsp?cboDataYear=2021&buttonYearandDataType=Search')

('2020', 'https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_results.jsp?cboDataYear=2020&buttonYearandDataType=Search')
```

**Step 3: Get the list (i.e., year_url_list) of tuples to obtain "HC-" zip file links**

**Lines 38-39**: The `base_hc_url` variable remains the same, representing the base URL for constructing Zip files' download URLs. The `zip_link_list` is an empty list that will store dictionaries containing information about the ZIP file download links.

**Lines 40-43**: The code then enters a `for` loop that iterates through each tuple in `year_url_list`. Inside the loop, it sends a GET request to the year-specific URL, and if the request is successful (`response.raise_for_status()`), it proceeds to extract the data using `pd.read_html()`.

**Lines 45-46**: The `pd.read_html()` function parses the HTML response text and attempts to extract tables from it. Each table is represented as a DataFrame in the `dfs` list. Next, the code initializes an empty DataFrame called `hc_df` to store the relevant data from the tables.

**Lines 47-53**: For each DataFrame `df` in `dfs`, the code checks if it contains the column "File(s), Documentation & Codebooks". If it does, the code concatenates `df` with `hc_df`, drops rows with missing values in the "PUF no." column, selects only rows where the "PUF no." contains "HC-", drops duplicate rows based on the "PUF no.", and appends the MEPS PUF URL by combining `base_hc_url` with the values in the "PUF no." column.

**Lines 54-66**: If `hc_df` is not empty, the code renames the columns of `hc_df` and iterates over each row using `hc_df.itertuples()`. For each row, it sends a GET request to the healthcare data file URL (`row.hc_url`), raises an exception if the request fails, and parses the HTML response using `BeautifulSoup`.

**Lines 68-83:** If successful, the code extracts the name of the MEPS-HC PUF file from the HTML using `hc_soup.find(class_="OrangeBox").text` and searches for a specific table cell ("td") containing the text "Data File". If found, it creates a dictionary called `zip_link_dict` and populates it with information such as the PUF number, data year, PUF name, file format, and the ZIP file download link. The `zip_link_dict` is then appended to `zip_link_list`.

The code below prints selected elements from `zip_link_list`.

```
for item in zip_link_list[:2]:   print(item)
```

The output:

```
{'PUF_num': 'HC-228', 'data_year': '2021', 'meps_file': 'MEPS HC-228: 2021 Full Year Popul
ation Characteristics File', 'file_format': 'Data File, ASCII format', 'zip_link': 'https:
//meps.ahrq.gov/data_files/pufs/h228/h228dat.zip'}

{'PUF_num': 'HC-228', 'data_year': '2021', 'meps_file': 'MEPS HC-228: 2021 Full Year Popul
ation Characteristics File', 'file_format': 'Data File, SAS transport format', 'zip_link':
 'https://meps.ahrq.gov/data_files/pufs/h228/h228ssp.zip'}
```

**Lines 84-90:** Any exceptions that occur during the execution of the code are caught using `try`/`except` blocks and error messages are printed.

Overall, this code retrieves ZIP files' download URLs for each available year based on the information obtained from the previous code and appends the relevant information to the `zip_link_list`.

**Step 4: Create an Excel file from the Pandas DataFrame**

**Line 100:** The code first creates the DataFrame `meps_df` using the `pd.DataFrame()` constructor and passing `zip_link_list` as an argument, which contains the information about the ZIP file download links.

**Line 101:** The code removes duplicate rows in `meps_df` using the `drop_duplicates()` method with the `inplace=True` argument to modify the DataFrame in-place.

**Line 102:** The code modifies the "meps_file" column in `meps_df` by splitting the values using `str.split(", ", n=1)` and selecting the second part of the split using `str[-1]`. It extracts the MEPS file name and updates the values in the "meps_file" column.

9

**Line 103:** The code modifies the "file_format" column in `meps_df` by splitting the values using `str.split(", ", n=1)` and selecting the second part of the split using `str[-1]`. It extracts the file format and updates the values in the "file_format" column.

**Line 104**: The code performs multi-key sorting using the 'sort_values' attribute with the ascending= argument using True or False for the same number of values.

**Line 105:** The code prints the first few rows of the modified DataFrame using the `head()` method.

```
data_year puf_num      meps_file          file_format  \
20      2021  HC-227  2021 Jobs File          ASCII format
22      2021  HC-227  2021 Jobs File         SAS V9 format
21      2021  HC-227  2021 Jobs File  SAS transport format
23      2021  HC-227  2021 Jobs File          Stata format
24      2021  HC-227  2021 Jobs File           XLSX format


                                            zip_link
20  https://meps.ahrq.gov/data_files/pufs/h227/h22...
22  https://meps.ahrq.gov/data_files/pufs/h227/h22...
21  https://meps.ahrq.gov/data_files/pufs/h227/h22...
23  https://meps.ahrq.gov/data_files/pufs/h227/h22...
24  https://meps.ahrq.gov/data_files/pufs/h227/h22...
```

**Line 109:** The variable `today` represents the current date. It uses `pd.Timestamp("now")` to obtain the current timestamp and `strftime("%Y-%m-%d")` to format it as "YYYY-MM-DD".

**Line 110**: The `to_excel()` method is used to write the DataFrame to an Excel file, with the `index=False` argument to exclude the DataFrame index from the output. The code saves the DataFrame to an Excel file with a dynamically generated filename. It uses f-string formatting to include the `today` variable in the filename, resulting in a filename like "C:\SESUG_2013\Python_Solution_WS_YYYY-MM-DD.xlsx".

Overall, this code modifies and displays the DataFrame `meps_df`, generates the current date, and saves the DataFrame to an Excel file with a filename that includes the current date in the "output" folder.

## Comparing SAS and Python Output

We read both SAS and Python-generated Excel files using read_excel() method of pandas into pandas DataFrames. We then compare the two DataFrames using pandas compare(). This comparison confirms that SAS and Python applications' file download URL lists are identical.

```
#!/usr/bin/env python
# coding: utf-8

import pandas as pd
df_sas = pd.read_excel(f"C:\SESUG_2023\SAS_Solution_WS_2023-06-08.xlsx")
df_py = pd.read_excel(f"C:\SESUG_2023\Python_Solution_WS_2023-06-08.xlsx")
diff = df_sas.compare(df_py, keep_equal=True, keep_shape = True)
print(diff)
```

The partial output:

```
      data_year  ...                                          zip_link
            self  ...                                            other
0            2021  ...  https://meps.ahrq.gov/data_files/pufs/h227/h22...
1            2021  ...  https://meps.ahrq.gov/data_files/pufs/h227/h22...
2            2021  ...  https://meps.ahrq.gov/data_files/pufs/h227/h22...
3            2021  ...  https://meps.ahrq.gov/data_files/pufs/h227/h22...
4            2021  ...  https://meps.ahrq.gov/data_files/pufs/h227/h22...
      ...    ...                                                  ...
1174         1996  ...    https://meps.ahrq.gov/data_files/pufs/h17ssp.zip
1175         1996  ...    https://meps.ahrq.gov/data_files/pufs/h24dat.zip
1176         1996  ...    https://meps.ahrq.gov/data_files/pufs/h24ssp.zip
1177         1996  ...    https://meps.ahrq.gov/data_files/pufs/h41dat.zip
1178         1996  ...    https://meps.ahrq.gov/data_files/pufs/h41ssp.zip

[1179 rows x 10 columns]
```

## CONCLUSION

The paper offers SAS® and Python solutions to dynamically create a list of MEPS-HC files' download URLs representing various formats (e.g., ASCII, Excel, SAS transport, SAS V9, and Stata files). This comparison confirms that SAS and Python applications' file download URL lists are identical This work found that code efficiencies varied between the two programs. For example, Python's libraries BeautifulSoup and urllib.request for HTTP Requests came in handy to perform web scraping. On the other hand, in performing the same task, some users may find the SAS program much more intuitive. However, we see the Python solution as the most efficient due to its flexibility, ease of use, and minimal coding.

These programs are reusable to obtain an updated list of MEPS-HC ZIP files' download links as more data sets are added to the website annually from the ongoing panel survey. In other words, future use of our programs will likely result in more such URLs depending on the timing of their execution, as the target web page and the linked websites get updated with new files. Once you automate the creation of a list of all MEPS-HC ZIP files' download URLs in an Excel file using one of these programs (SAS or Python), you can directly initiate ZIP file downloads by simply clicking the relevant URLs from that file outside of the Internet, an easy alternative to navigating multiple pertinent websites for the same task.

## REFERENCES

1. Cohen, J., S. Cohen, and J. Banthin. 2009. "The Medical Expenditure Panel Survey: A National Information Resource to Support Health Care Cost Research and Inform Policy and Practice." *Medical Care*, 47:S44-50.

## ACKNOWLEDGMENTS

## CONTACT INFORMATION

Your comments are valued and encouraged. Contact the authors at:

Pradip K. Muhuri, PhD
AHRQ/CFACT
5600 Fishers Lane
Rockville, MD 20857
Pradip.Muhuri@ahrq.hhs.gov


Charles Z.X. Han
Ryman Healthcare Ltd.
charleszhouxiaohan@gmail.com


John Vickery
North Carolina State University Libraries
(919) 513-0344
John_vickery@ncsu.edu

```
/* APPENDIX I (SAS Program) */
1  %let path = c:\SESUG_2023;
2  /* Step 1: Fetch the main MEPS-HC web page's contents using PROC HTTP */
3  filename source "&path/web_file.txt";
4  proc http
5      url = "https://meps.ahrq.gov/data_stats/download_data_files.jsp"
6      out=source;
7  run;
8
9  /* Step 2:  Parse the MEPS-HC main web page */
10 /* Get years from dropdown options list */
11 data year_values;
12     length year $4;
13     infile source length = reclen lrecl = 32767 end=eof;
14     re = prxparse('/<option value="\d{4}">(\d{4})<\/option>/');
15     /* Read the HTML line by line */
16     do while (not eof);
17             input html_line $varying32767. reclen;
18         /* Match and extract the years using regular expressions */
19         if prxmatch(re, html_line) > 0 then do;
20             call prxposn(re, 1, start, end);
21             year = substr(html_line, start, end);
22             output year_values;
23         end;
24     end;
25 run;
26
27 /* De-dup and keep only the year variable in the data set */
28 proc sort data=year_values (keep=year) nodupkey;
29     by _ALL_;
30 run;
31
32 /* Step 3: Call Proc HTTP for each year_url */
33 %macro get_year(year);
34     filename yearresp "&path\hc_response.txt";
35     %local base_url data_year_param search_param year_url;
36     /* URL elements */
37     %let base_url = https://meps.ahrq.gov/data_stats/download_data_files_results.jsp?;
38     %let data_year_param = cboDataYear=&year.;
39     %let search_param = %nrstr(&buttonYearandDataType=Search);
40
41     /* Combine the URL elements to construct the year url*/
42     %let year_url = &base_url.&data_year_param.&search_param.;
43
44     /* Call PROC HTTP to retrieve the content */
45     /* of each year search results */
46     proc http
47             url = "&year_url."
48     out=yearresp;
49     run;
50
51     data year_list_&year. (keep=puf_num meps_file data_year);
52             length puf_num $15 meps_file $150 data_year $10;
53             infile yearresp length = reclen lrecl = 32767 end=eof;
54             /* regex to get the PUF num in the table of results */
55             prx_puf = prxparse('/<a href="download_data_files_detail\.jsp\?cboPufNumber=.+\">(.+)<\/a>/');
56             /* regex to get the meps file */
```

```
57              prx_meps = prxparse('/<td height="30" align="left" valign="top"
class="bottomRightgrayBorder"><div align="left" class="contentStyle">(.+)<\/font>/');
58                   /* regex for data year */
59                   prx_data_year = prxparse('/<td width="1" height="30" align="left" valign="top"
class="bottomRightgrayBorder"><div align="left" class="contentStyle">(.+)<\/font><\/div><\/td>/');
60                   /* Read the HTML line by line */
61                   do while (not eof);
62                        input html_line $varying32767. reclen;
63                        if prxmatch(prx_puf, html_line) > 0 then do;
64                             call prxposn(prx_puf, 1, start, end);
65                        puf_num = substr(html_line, start, end);
66                        end;
67                     if prxmatch(prx_meps, html_line) > 0 then do;
68                        call prxposn(prx_meps, 1, start, end);
69                        meps_file = substr(html_line, start, end);
70                        meps_file = prxchange('s/<.+>//', -1, meps_file);        /* remove any stray html tags */
71                     end;
72                     if prxmatch(prx_data_year, html_line) > 0 then do;
73                             call prxposn(prx_data_year, 1, start, end);
74                        data_year = substr(html_line, start, end);
75                        output year_list_&year.;
76                        end;
77                end;
78      run;
79
80      filename yearresp clear;
81
82 %mend get_year;
83
84 /* Loop through each year of the year_values dataset and call the macro */
85 data _null_;
86      set year_values;
87    call execute('%nrstr(%get_year('||strip(year)||'));');
88 run;
89
90 /* Concatenate all year_list_YEAR datasets */
91 /* Just keep obs with puf_num beginning HC- */
92 data year_list;
93      set year_list_: ;
94      if substr(puf_num, 1, 2) = 'HC';
95 run;
96
97 /* De-dup */
98 proc sort data=year_list nodupkey;
99   by _ALL_;
100  run;
101
102 /* Step 4: Get the zip files for each HC- puf_num. */
103 %macro get_puf(pufnum);
104    %local base_url puf_url;
105    filename pufresp "&path\puf_response.txt";
106
107      /* URL elements */
108      %let base_url =
https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber=;
109
110      /* Combine the URL elements to construct the PUF url */
111      %let puf_url = &base_url.&pufnum.;
```

```
112
113     proc http
114             url="&puf_url."
115             out= pufresp;
116     run;
117
118     data puf_list_%sysfunc(translate(&pufnum, '_', '-')) (keep=puf_num file_format zip_link);
119             length puf_num $15 file_format $150 zip_link $200;
120             infile pufresp length = reclen lrecl = 32767 end=eof;
121
122             prx_data_file = prxparse('/<td width="50%" height="0" class="bottomRightgrayBorder">Data File.*,
(.+)<\/td>/');
123             prx_zip = prxparse('/<a href="\.\.\/(data_files[\/pufs]*\/.+\.zip)">ZIP<\/a>/');
124
125             do while (not eof);
126                     input html_line $varying32767. reclen;
127                     puf_num = "&pufnum.";
128                     if prxmatch(prx_data_file, html_line) > 0 then do;
129                             call prxposn(prx_data_file, 1, start, end);
130                     file_format = substr(html_line, start, end);
131                     end;
132
133                     if prxmatch(prx_zip, html_line) > 0 then do;
134                             call prxposn(prx_zip, 1, start, end);
135                     zip_link = cats("https://meps.ahrq.gov/", substr(html_line, start, end));
136                     output puf_list_%sysfunc(translate(&pufnum, '_', '-'));
137                     end;
138             end;
139
140     run;
141
142     filename pufresp clear;
143
144 %mend get_puf;
145
146 /* Loop through each puf_num of the year_list dataset and call the macro */
147 data _null_;
148     set year_list;
149     call execute('%nrstr(%get_puf('||strip(puf_num)||'));');
150 run;
151
152 /* Concatenate all puf_list datasets */
153 data puf_list;
154     set puf_list_: ;
155 run;
156
157 /* De-dup */
158 proc sort data=puf_list nodupkey;
159     by _ALL_;
160 run;
161
162 /* Step 5: Merge year_list and puf_list */
163 proc sql;
164     create table meps_zip_links as
165     select y.puf_num, y.meps_file, y.data_year, p.file_format, p.zip_link
166     from year_list as y,
167             puf_list as p
168     where y.puf_num = p.puf_num
```

```
169    order by data_year desc, puf_num, file_format;
170  quit;
171
172  /* Format the current date as "YYYY-MM-DD" */
173  %let current_date = %sysfunc(today());
174  %let formatted_date = %sysfunc(putn(&current_date., yymmdd10.));
175
176  /* Step 6: Direct proc report output to excel */
177  ods listing close;
178  ods excel file = "&path\SAS_Solution_WS_&formatted_date..xlsx"
179    options (sheet_name = 'Sheet1'
180    flow="header,data" row_heights = '15'
181    absolute_column_width='11,11,70,30,55');
182  proc report data=meps_zip_links;
183    column data_year puf_num meps_file file_format  zip_link;
184    define puf_num / display ;
185    define meps_file / display;
186    define data_year / display;
187    define file_format / display;
188    define zip_link / display ;
189    compute zip_link ;
190      call define(_col_,"url",zip_link);
191    endcomp;
192  run;
193  ods excel close;
194  ods listing;
```

**APPENDIX II (Part of the Excel file from the PROC REPORT Output)**

| data_year | PUF_num | meps_file | file_format | zip_link |
|---|---|---|---|---|
| 2021 | HC-228 | Full Year Population Characteristics File | ASCII format | https://meps.ahrq.gov/data_files/pufs/h228/h228dat.zip |
| 2021 | HC-228 | Full Year Population Characteristics File | SAS V9 format | https://meps.ahrq.gov/data_files/pufs/h228/h228v9.zip |
| 2021 | HC-228 | Full Year Population Characteristics File | SAS transport format | https://meps.ahrq.gov/data_files/pufs/h228/h228ssp.zip |
| 2021 | HC-228 | Full Year Population Characteristics File | Stata format | https://meps.ahrq.gov/data_files/pufs/h228/h228dta.zip |
| 2021 | HC-228 | Full Year Population Characteristics File | XLSX format | https://meps.ahrq.gov/data_files/pufs/h228/h228xlsx.zip |
| 2021 | HC-227 | Jobs File | ASCII format | https://meps.ahrq.gov/data_files/pufs/h227/h227dat.zip |
| 2021 | HC-227 | Jobs File | SAS V9 format | https://meps.ahrq.gov/data_files/pufs/h227/h227v9.zip |
| 2021 | HC-227 | Jobs File | SAS transport format | https://meps.ahrq.gov/data_files/pufs/h227/h227ssp.zip |
| 2021 | HC-227 | Jobs File | Stata format | https://meps.ahrq.gov/data_files/pufs/h227/h227dta.zip |
| 2021 | HC-227 | Jobs File | XLSX format | https://meps.ahrq.gov/data_files/pufs/h227/h227xlsx.zip |

**Appendix III (Python program)**

```python
1   #!/usr/bin/env python
2   # coding: utf-8
3
4   # **step 1**
5   # - note the python libraries/modules
6
7   # import libraries
8   import pandas as pd
9   import requests
10  from bs4 import BeautifulSoup
11  from urllib.parse import quote
12
13  # **step 2**
14  # - get the list of numerical year options from the main page
15  # - save as list of tuples with year and url
16
17  # main page response
18  main_page = requests.get("https://meps.ahrq.gov/data_stats/download_data_files.jsp")
19  main_soup = BeautifulSoup(main_page.text, "html.parser")
20
21  # get list of data years and links for each year
22  # skip the "All years" (i.e. non-digit options)
23  base_year_url =
    "https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_results.jsp?cboDataYear="
24  year_url_suffix = "&buttonYearandDataType=Search"
25
26  year_options = main_soup.select_one('select[name="cboDataYear"]').select("option")
27
28  year_url_list = [
29      (x.get("value"), base_year_url + quote(x.get("value")) + year_url_suffix)
30      for x in year_options
31      if x.get("value").isdigit()
32  ]
33
34  # **step 3**
35  # - use the year_url_list list of tuples to get "HC-" zip file links
36  # - save these in list of dictionaries
37
38  base_hc_url =
    "https://meps.ahrq.gov/mepsweb/data_stats/download_data_files_detail.jsp?cboPufNumber="
39  zip_link_list = []
40  for year, year_url in year_url_list:
41      try:
42          response = requests.get(year_url)
43          response.raise_for_status()
44          try:
45              dfs = pd.read_html(response.text)
46              hc_df = pd.DataFrame()
47              for df in dfs:
48                  if "File(s), Documentation & Codebooks" in df.columns:
49                      hc_df = pd.concat([hc_df, df], ignore_index=True)
50                      hc_df = hc_df.dropna(subset="PUF no.")
```

18

```
51            hc_df = hc_df.loc[hc_df["PUF no."].str.contains("HC-")]
52            hc_df = hc_df.drop_duplicates(subset="PUF no.")
53            hc_df["hc_url"] = base_hc_url + hc_df["PUF no."]
54       if not hc_df.empty:
55          hc_df.columns = [
56             "puf_num",
57             "Files",
58             "Data_Update",
59             "Year",
60             "File_Type",
61             "hc_url",
62          ]
63          for row in hc_df.itertuples():
64             hc_response = requests.get(row.hc_url)
65             hc_response.raise_for_status()
66             hc_soup = BeautifulSoup(hc_response.text)
67             try:
68                meps_file = hc_soup.find(class_="OrangeBox").text
69                for td in hc_soup.find_all("td"):
70                   zip_link_dict = {}
71                   if td.text.startswith("Data File"):
72                      zip_link_dict["data_year"] = row.Year
73                      zip_link_dict["puf_num"] = row.puf_num
74                      zip_link_dict["meps_file"] = meps_file
75                      zip_link_dict["file_format"] = td.text
76                      zip_link_dict[
77                         "zip_link"
78                      ] = "https://meps.ahrq.gov" + td.find_next("a").get(
79                         "href"
80                      ).strip(
81                         ".."
82                      )
83                      zip_link_list.append(zip_link_dict)
84             except:
85                # catch your exceptions here
86                pass
87       except requests.exceptions.HTTPError as err:
88          print(err)
89    except requests.exceptions.HTTPError as httperr:
90       print(httperr)
91
92  # **step 4**
93  # - save the list of dictionaries to pandas dataframe
94  # - deduplicate
95  # - clean up meps_file column to remove "PUF no." prefix
96  # - clean up file_format column to remove "Data File" prefix
97  # - sort the dataframe in certain order
98  # - get n rows of the dataframe
99
100  meps_df = pd.DataFrame(zip_link_list)
101  meps_df.drop_duplicates(inplace=True)
102  meps_df["meps_file"] = meps_df["meps_file"].str.split(": ", n=1).str[-1]
103  meps_df["file_format"] = meps_df["file_format"].str.split(", ", n=1).str[-1]
104  meps_df = meps_df.sort_values(by=['data_year','puf_num', 'file_format'], ascending=[False, True,
True])
105  print(meps_df.head(5))
106
107  # **step 4 (continued)**
```

```
108  # - obtain the current date (timestamp)
109  # - save the dataframe as an excel file
110
111  today = pd.Timestamp("now").strftime("%Y-%m-%d")
112  meps_df.to_excel(f"C:\SESUG_2023\Python_Solution_WS_{today}.xlsx", index=False)
```