```python
  1 #!/usr/bin/env python3
  2 # -*- coding: utf-8 -*-
  3 """
  4 Analyze the structure and broad properties of PCD online archive, including
  5 Spanish-language, some of which are indexed separately
  6
  7 @author: chadheilig
  8
  9 Sections of this script, based on levels of MMWR archive:
 10 0. PCD home https://www.cdc.gov/pcd/index.htm
 11 1. Contents of series, including Spanish; list of archive volumes
 12 2. List and contents of volumes
 13 3. List of articles
 14 4. Complete list of PCD files
 15
 16 Main product: pcd_dframe
 17 """
 18
 19 #%% Import modules and set up environment
 20 # import from 0_cdc-corpora-header.py
 21
 22 os.chdir('/Users/cmheilig/cdc-corpora/_test')
 23
 24 #%% 0. Start with PCD home https://www.cdc.gov/pcd/index.htm
 25 base_url = 'https://www.cdc.gov/pcd/index.htm'
 26 home_a = BeautifulSoup(get_html_from_url(base_url), 'lxml').\
 27     find('a', href=re.compile('pcd/index.htm'),
 28             string=re.compile('Preventing Chronic Disease'))
 29 # process_aTag(home_a, base_url)
 30 # {'base': 'https://www.cdc.gov/pcd/index.htm',
 31 #  'href': '/pcd/index.htm',
 32 #  'url': 'https://www.cdc.gov/pcd/index.htm',
 33 #  'path': '/pcd/index.htm',
 34 #  'filename': 'index.htm',
 35 #  'mirror_path': '/pcd/index.htm',
 36 #  'string': 'Preventing Chronic Disease'}
 37
 38 home_dframe = pd.DataFrame(process_aTag(home_a, base_url), index = [0])
 39 # home_dframe.loc[:, ['path', 'string']]
 40 #                path                        string
 41 # 0  /pcd/index.htm  Preventing Chronic Disease
 42 home_html = get_html_from_url(home_dframe.url[0]) # len(home_html) # 188351
 43 home_soup = BeautifulSoup(home_html, 'lxml')
 44
 45 # review all anchor-hrefs from home URL
 46 # len(home_soup.find_all('a', href=True)) # 110
 47 # pd.DataFrame([process_aTag(aTag, home_dframe.url[0])
 48 #     for aTag in home_soup.find_all('a', href=True)]).\
 49 #     to_excel('pcd-home-anchors.xlsx', engine='openpyxl')
 50 # [110 rows x 7 columns]
 51
 52 #%% 1. Contents of series, including Spanish; list of archive volumes
 53
 54 # Review of anchor elements in home page, pcd-home-anchors.xlsx
 55 # https://www.cdc.gov/pcd/current_issue.htm  # current volume
 56 #     all issues and articles in 2020 (to date)
```

```python
 57  # https://www.cdc.gov/pcd/issues/archive.htm # past volumes
 58  #    all volumes and articles in 2004-2011, volumes in 2012-2019
 59
 60  series_a = home_soup.find_all('a', href=re.compile('archive'))
 61  # [<a href="/pcd/issues/archive.htm">Issue Archive</a>]
 62
 63  # Home page does not point to Spanish-language archive
 64  # https://www.cdc.gov/pcd/es/archive_es.htm
 65  home_es_url = 'https://www.cdc.gov/pcd/es/archive_es.htm'
 66  series_es_a = BeautifulSoup(get_html_from_url(home_es_url), 'lxml').\
 67      find('a', href=re.compile('pcd/es/archive_es.htm'))
 68  # process_aTag(series_es_a, home_es_url)
 69  # {'base': 'https://www.cdc.gov/pcd/es/archive_es.htm',
 70  #  'href': '/pcd/es/archive_es.htm',
 71  #  'url': 'https://www.cdc.gov/pcd/es/archive_es.htm',
 72  #  'path': '/pcd/es/archive_es.htm',
 73  #  'filename': 'archive_es.htm',
 74  #  'mirror_path': '/pcd/es/archive_es.htm',
 75  #  'string': 'Archivo de números en español'}
 76
 77  series_dframe = pd.DataFrame(
 78      [process_aTag(series_a[0], home_dframe.url[0]),
 79       process_aTag(series_es_a, home_es_url)])
 80  # series_dframe.loc[:, ['path', 'string']]
 81  #                         path                        string
 82  # 0  /pcd/issues/archive.htm                Issue Archive
 83  # 1   /pcd/es/archive_es.htm  Archivo de números en español
 84
 85  series_html = [get_html_from_url(url) for url in series_dframe.url]
 86  # [len(x) for x in series_html] # [210632, 35870]
 87  series_soup = [BeautifulSoup(html, 'lxml') for html in series_html]
 88
 89  # review all anchor-hrefs from series URL
 90  # pd.DataFrame([process_aTag(aTag, url)
 91  #      for soup, url in zip(series_soup, series_dframe.url)
 92  #      for aTag in soup.find_all('a', href=True)]).\
 93  #      to_excel('pcd-series-anchors.xlsx', engine='openpyxl')
 94  # [316 rows x 7 columns]
 95
 96  #%% 2. List and contents of volumes
 97
 98  # Review of anchor elements in series page, pcd-series-anchors.xlsx
 99  # https://www.cdc.gov/pcd/current_issue.htm  # current volume
100  #    current volume, all issues in 2021 (to date)
101  # https://www.cdc.gov/pcd/issues/yyyy/yyyy_TOC.htm
102  #    all volumes and articles in 2012-2020
103  # https://www.cdc.gov/pcd/issues/yyyy/mmm/toc.htm
104  #    all volumes and articles in 2004-2011
105  # https://www.cdc.gov/pcd/es/yyyy_toc.htm
106  #    all volumes and articles in 2012-2014
107  # https://www.cdc.gov/pcd/es/yyyy_mmm_toc.htm
108  #    all volumes and articles in 2005-2011
109  # https://www.cdc.gov/pcd/spanish/current_issue_es.htm # can ignore
110  #    last updated 2015
111
112  pcd_vol_re = re.compile(r'(current_issue|\d{4}.*(TOC|toc)).htm')
```

```
113  volumes_a = [soup.find_all('a', href=pcd_vol_re) for soup in series_soup]
114  volumes_a_n = [len(x) for x in volumes_a] # sum(volumes_a_n) # [47, 36]
115
116  volumes_dframe = pd.DataFrame([process_aTag(aTag, url)
117      for a_list, url in zip(volumes_a, series_dframe.url)
118      for aTag in a_list])
119  # volumes_dframe.loc[:, ['path', 'string']]
120  # with pd.option_context("display.max_rows", 100):
121  #     display(volumes_dframe.loc[:, ['path', 'string']])
122  #                             path              string
123  # 0    /pcd/issues/2020/2020_TOC.htm              2020
124  # ..                            ...               ...
125  # 8    /pcd/issues/2012/2012_TOC.htm              2012
126  # 9     /pcd/issues/2011/nov/toc.htm          November
127  # ..                            ...               ...
128  # 45    /pcd/issues/2004/jan/toc.htm           January
129  # 46            /pcd/current_issue.htm  View Current Volume
130  # 47             /pcd/es/2014_toc.htm              2014
131  # 48             /pcd/es/2013_toc.htm              2013
132  # 49             /pcd/es/2012_toc.htm              2012
133  # 50         /pcd/es/2011_nov_toc.htm         Noviembre
134  # ..                            ...               ...
135  # 82         /pcd/es/2005_jan_toc.htm             Enero
136
137  volumes_html = [get_html_from_url(url) for url in tqdm(volumes_dframe.url)]
138  # 83/83 [00:15<00:00,  5.46it/s]
139  # [len(x) for x in volumes_html]
140  # [336609, 328913, 407823, 393428, 424664, 466946, 191632, 176170, 148415, ...]
141  volumes_soup = [BeautifulSoup(html, 'lxml') for html in volumes_html]
142
143  # review all anchor-refs from volumes URLs
144  # pd.DataFrame([process_aTag(aTag, url)
145  #     for soup, url in zip(volumes_soup, volumes_dframe.url)
146  #     for aTag in soup.find_all('a', href=True)]).\
147  #     to_excel('pcd-volumes-anchors.xlsx', engine='openpyxl')
148  # [10687 rows x 7 columns]
149
150  #%% 3. List of articles
151
152  # Review of anchor elements in volumes page, pcd-volumes-anchors.xlsx
153  # mostly filenames of form dd_dddd.htm or dd_dddd_es.htm
154  # Retrieve files under https://www.cdc.gov/pcd/issues/
155  # Regular expressions for full paths
156  #     \d{4}/(jan|mar|apr|may|jul|sep|oct|nov)/ # 2004-2011
157  #     \d{4}/                                   # 2012-2021
158  #         \d{2}_\d{4,5}([aber]|_es)?.htm
159  # Regular expression for hrefs: \d{2}_\d{4,5}([aber]|_es)?.htm
160  # These include Spanish but exclude French (357), Portuguese (1),
161  #     Vietnamese (1), and Chinese (simplified [356] and traditional [356]),
162  #         \d{2}_\d{4}_(fr|pr|vi|zhs|zht).htm
163  # _es last seen in 2014; other language suffixes last seen Jan 2010
164
165  pcd_art_re = re.compile(r'\d{2}_\d{4,5}([aber]|_es)?.htm')
166  articles_a = [soup.find_all('a', href=pcd_art_re) for soup in volumes_soup]
167  articles_a_n = [len(x) for x in articles_a] # sum(articles_a_n) # 4405
168  # [170, 166, 166, 142, 181, 231, 230, 216, 179, 67, 40, 38, 52, 42, 53, ...]
```

```
169
170 articles_dframe = pd.DataFrame([process_aTag(aTag, url)
171    for a_list, url in zip(articles_a, volumes_dframe.url)
172    for aTag in a_list])
173 # (4405, 7)
174 # articles_dframe.loc[:, ['path', 'string']]
175 with pd.option_context("display.max_colwidth", 36):
176     display(articles_dframe.loc[:, ['path', 'string']])
177 #                                   path                      string
178 # 0              /pcd/issues/2020/20_0214.htm  Collecting Early Childhood Obesi...
179 # 1              /pcd/issues/2020/20_0262.htm  Perceived Importance of Physical...
180 # 2              /pcd/issues/2020/19_0431.htm  Chronic Disease Among African Am...
181 # 3              /pcd/issues/2020/20_0366.htm  Water Safety in California Publi...
182 # 4              /pcd/issues/2020/20_0340.htm  "We're, Like, the Most Unhealthy...
183 #                                    ...                      ...
184 # 4400  /pcd/issues/2005/jan/04_0079_es.htm  De la investigación a la práctic...
185 # 4401  /pcd/issues/2005/jan/04_0075_es.htm  Pasos Adelante:|La eficacia de u...
186 # 4402  /pcd/issues/2005/jan/04_0076_es.htm  El índice de sanidad escolar (Sc...
187 # 4403  /pcd/issues/2005/jan/04_0083_es.htm  El desarrollo y la adaptación de...
188 # 4404  /pcd/issues/2005/jan/04_0077_es.htm  La|Border Health Strategic Initi...
189
190 # Check for duplicate URLs
191 articles_repeated = {
192    label: content.loc[content.duplicated(keep = False)].index.to_list()
193       for label, content
194       in articles_dframe.loc[:, ['href', 'url', 'path', 'filename']].items() }
195 # { k: len(v) for k, v in articles_repeated.items() }
196 # {'href': 26, 'url': 1429, 'path': 1429, 'filename': 1443}
197
198 # 4 articles (10 records) have same referring source and same target
199 dupes = articles_dframe.duplicated(['base', 'url'], keep=False) # dupes.sum() # 10
200
201 # 2008/jan/06_0177, 2008/jan/06_0177_es, 2018/17_0395, 2020/19_0176
202 with pd.option_context("display.max_colwidth", 65):
203     display(articles_dframe.loc[dupes, 'url'])
204 # 163              https://www.cdc.gov/pcd/issues/2020/19_0176.htm
205 # 164              https://www.cdc.gov/pcd/issues/2020/19_0176.htm
206 # 454              https://www.cdc.gov/pcd/issues/2018/17_0395.htm
207 # 468              https://www.cdc.gov/pcd/issues/2018/17_0395.htm
208 # 2690         https://www.cdc.gov/pcd/issues/2008/jan/06_0177.htm
209 # 2691         https://www.cdc.gov/pcd/issues/2008/jan/06_0177.htm
210 # 2693         https://www.cdc.gov/pcd/issues/2008/jan/06_0177.htm
211 # 4154      https://www.cdc.gov/pcd/issues/2008/jan/06_0177_es.htm
212 # 4155      https://www.cdc.gov/pcd/issues/2008/jan/06_0177_es.htm
213 # 4156      https://www.cdc.gov/pcd/issues/2008/jan/06_0177_es.htm
214
215 # on review:
216 # 2008/jan/06_0177: drop 2690, 2693 as erroneous (no anchor text); keep 2691
217 # 2008/jan/06_0177_es: drop 4154, 4156 as erroneous (no anchor text); keep 4155
218 # 2018/17_0395: drop 468, keep 454 (duplicate references to same article)
219 # 2020/19_0176: drop 164, keep 163 (duplicate references to same article)
220
221 articles_dframe.drop([2690, 2693, 4154, 4156, 468, 164], inplace=True)
222 # (4399, 7)
223 articles_dframe.duplicated(['base', 'url'], keep=False).sum() # 0
224
```

```
225  # Further review reveals that every URL referenced from the Spanish-language
226  # archive should end in _es.htm rather than just .htm. 16 URLs are incorrect.
227  from_es_src = articles_dframe.base.str.contains('/es/') # 1011 True, 3388 False
228  has_es_name = articles_dframe.path.str.contains('_es.htm') # 1703 True, 2696 False
229  pd.crosstab(from_es_src, has_es_name, margins=True).iloc[[1,0,2],[1,0,2]]
230  # path    True   False    All
231  # base
232  # True     995      16   1011
233  # False    708    2680   3388
234  # All     1703    2696   4399
235
236  # 16 referenced from Spanish-language archive, not named *_es.htm
237  es_src_not_name = from_es_src & ~has_es_name
238  articles_dframe.loc[es_src_not_name, ['base', 'filename']] # mostly 2005/jul
239  #                                                base       filename
240  # 3571          https://www.cdc.gov/pcd/es/2012_toc.htm  11_0345.htm
241  # 3572          https://www.cdc.gov/pcd/es/2012_toc.htm  12_0010.htm
242  # 4360  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0021.htm
243  # 4361  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0146.htm
244  # 4362  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0127.htm
245  # 4363  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0144.htm
246  # 4364  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0136.htm
247  # 4365  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0130.htm
248  # 4366  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0124.htm
249  # 4367  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0009.htm
250  # 4368  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0129.htm
251  # 4369  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0137.htm
252  # 4370  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0126.htm
253  # 4371  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0003.htm
254  # 4372  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0023.htm
255  # 4373  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0121.htm
256
257  # revise url, path, filename, mirror_path
258  articles_dframe.href.loc[es_src_not_name] = \
259     articles_dframe.href.loc[es_src_not_name].str.replace('.htm', '_es.htm')
260  articles_dframe.url.loc[es_src_not_name] = \
261     articles_dframe.url.loc[es_src_not_name].str.replace('.htm', '_es.htm')
262  articles_dframe.path.loc[es_src_not_name] = \
263     articles_dframe.path.loc[es_src_not_name].str.replace('.htm', '_es.htm')
264  articles_dframe.filename.loc[es_src_not_name] = \
265     articles_dframe.filename.loc[es_src_not_name].str.replace('.htm', '_es.htm')
266  articles_dframe.mirror_path.loc[es_src_not_name] = \
267     articles_dframe.mirror_path.loc[es_src_not_name].str.replace('.htm', '_es.htm')
268
269  articles_dframe.loc[es_src_not_name, ['base', 'filename']]
270  #                                                base         filename
271  # 3571          https://www.cdc.gov/pcd/es/2012_toc.htm  11_0345_es.htm
272  # 3572          https://www.cdc.gov/pcd/es/2012_toc.htm  12_0010_es.htm
273  # 4360  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0021_es.htm
274  # 4361  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0146_es.htm
275  # 4362  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0127_es.htm
276  # 4363  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0144_es.htm
277  # 4364  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0136_es.htm
278  # 4365  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0130_es.htm
279  # 4366  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  04_0124_es.htm
280  # 4367  https://www.cdc.gov/pcd/es/2005_jul_toc.htm  05_0009_es.htm
```

```
281  # 4368   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   04_0129_es.htm
282  # 4369   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   04_0137_es.htm
283  # 4370   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   04_0126_es.htm
284  # 4371   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   05_0003_es.htm
285  # 4372   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   05_0023_es.htm
286  # 4373   https://www.cdc.gov/pcd/es/2005_jul_toc.htm   04_0121_es.htm
287
288  # articles_dframe.href[es_src_not_name]
289  with pd.option_context("display.max_colwidth", 65):
290      display(articles_dframe.url[es_src_not_name])
291  # articles_dframe.path[es_src_not_name]
292  # articles_dframe.mirror_path[es_src_not_name]
293
294  # update cross-tabulation
295  has_es_name = articles_dframe.path.str.contains('_es.htm') # 1719 True, 2680 False
296  pd.crosstab(from_es_src, has_es_name, margins=True).iloc[[1,0,2],[1,0,2]]
297  # path    True   False    All
298  # base
299  # True    1011       0   1011
300  # False    708    2680   3388
301  # All     1719    2680   4399
302
303  # 708 referenced from English-language archive with name *_es.htm
304  es_name_not_src = ~from_es_src & has_es_name
305  este_string = articles_dframe.loc[es_name_not_src, ['mirror_path', 'string']]
306  { este: este_string.string.to_list().count(este)
307    for este in sorted(set(este_string.string)) }
308  # {'Este artículo en español': 54, 'Este resumen en español': 654}
309
310  # These 708 targets are referenced from Spanish-language archive, as well
311  # shown using asymmetric set difference
312  set(articles_dframe.url[~from_es_src & has_es_name]).\
313      difference(articles_dframe.url[from_es_src & has_es_name])
314  # len(set(articles_dframe.url[~from_es_src & has_es_name]).\
315  #     symmetric_difference(articles_dframe.url[from_es_src & has_es_name])) # 303
316  # Split articles_dframe: 3542 unique targets, 708 English-to-Spanish referents
317
318  articles_en_es_dframe = articles_dframe.loc[es_name_not_src] # (708, 7)
319  articles_dframe.drop(articles_dframe.index[es_name_not_src], inplace=True) # (3691, 7)
320
321  # Check again for duplicate URLs
322  articles_repeated = {
323     label: content.loc[content.duplicated(keep = False)].index.to_list()
324         for label, content
325         in articles_dframe.loc[:, ['href', 'url', 'path', 'filename']].items() }
326  # { k: len(v) for k, v in articles_repeated.items() }
327  # {'href': 8, 'url': 0, 'path': 0, 'filename': 20}
328  articles_dframe.loc[articles_repeated['href'], ['base', 'href']]
329  articles_dframe.loc[articles_repeated['filename'], ['filename', 'path']]
330  # based on path, these are duplicate names for distinct files
331
332  articles_dframe.index = list(range(3691))
333  articles_en_es_dframe.index = list(range(708))
334
335  articles_dframe.to_excel('pcd-articles_dframe.xlsx', engine='openpyxl')
336  articles_en_es_dframe.to_excel('pcd-articles_en_es_dframe.xlsx', engine='openpyxl')
```

```
337
338 #%% 4. Complete list of PCD files
339 pcd_dframe = pd.concat([
340     home_dframe.assign(level='home'),
341     series_dframe.assign(level='series'),
342     volumes_dframe.assign(level='volume'),
343     articles_dframe.assign(level='article'),
344     articles_en_es_dframe.assign(level='en_es')],
345     axis = 0, ignore_index = True)
346 # (4485, 8)
347
348 # pickle
349 pickle.dump(pcd_dframe, open("pcd_dframe.pkl", "wb"))
350 # pcd_dframe_ = pickle.load(open("pcd_dframe.pkl", "rb"))
351 # pcd_dframe.equals(pcd_dframe_)
352
353 # Excel; coulad also use engine=
354 pcd_dframe.to_excel('pcd_dframe.xlsx', engine='openpyxl')
355 # Excelternatives
356 # pcd_dframe.to_excel('pcd_dframe.xlsx', engine='xlsxwriter') # pd default
357 # pcd_dframe.to_excel('pcd_dframe.xls', engine='xlwt')
```