```python
1    #!/usr/bin/env python3
2    # -*- coding: utf-8 -*-
3
4    ###################################################################
5    # title:    BBS for Independence - data analysis using NLP
6    # subtitle: HSA in the ABC of Computational Text Analysis
7    # author:   Josias Bruderer, Universität Luzern
8    # date:     30. August 2021
9    # desc:     this script manages the whole analisis of textfiles.com
10   ###################################################################
11
12   # Preparations [R1]
13   """
14   The following lines of code is used for preparing our environment.
15   """
16
17   # load the necessary libraries
18   import os
19   import shutil
20   import sys
21   import time
22   from pathlib import Path
23   import textacy
24   import spacy
25   import pandas as pd
26   import scattertext as st
27   from plotnine import *
28   import numpy as np
29   from multiprocessing import Pool
30   from multiprocessing.managers import BaseManager
31   import string
32   from octis.preprocessing.preprocessing import Preprocessing
33   # from octis.dataset.dataset import Dataset
34   # from octis.evaluation_metrics.diversity_metrics import TopicDiversity
35   # from octis.models.LDA import LDA
36   # from wordcloud import WordCloud
37   import csv
38
39   project_path = Path.cwd()
40
41   # prepare to load project specific libraries
42   if project_path not in sys.path:
43       sys.path.append(str(project_path))
44
45   # import modules
46   from modules import wrangler
47   from modules import helpers
48   from modules import nlp_pool
49   from modules import years
50
51   # Data Wrangling [R2]
52   """
53   In this section the required data is downloaded and preprocessed (f.E. unzipped).
54   The module data_wrangler will be used for this.
55   """
56
57   number_of_threads = 24
58
59   """
60   skip_steps = [] # skip nothing
61   skip_steps = ["download", "cleaning"] # use this after modification on metadata_file_filter
62   skip_steps = ["download", "cleaning", "metadata-filtering", "modeling", "analysis_freq", "analysis_advance_preparation",
63                 "analysis_scattertext", "analysis_year", "analysis_octis", "analysis_entities"] # the full list
64   """
65   skip_steps = [] # the full list
66
67   data_url = "http://archives.textfiles.com/[name].zip"
68   data_names = ["100", "adventure", "anarchy", "apple", "art", "artifacts", "bbs", "computers", "conspiracy", "digest",
69                 "drugs", "etext", "exhibits", "floppies", "food", "fun", "games", "groups", "hacking", "hamradio",
70                 "history", "holiday", "humor", "internet", "law", "magazines", "media", "messages", "music", "news",
71                 "occult", "phreak", "piracy", "politics", "programming", "reports", "rpg", "science", "sex", "sf",
72                 "stories", "survival", "tap", "ufo", "uploads", "virus",
73                 "fidonet-on-the-internet"]  # categories to download
74   data_names_exclude = ["fidonet-on-the-internet", "tap", "floppies", "exhibits", "artifacts",
75                         "piracy", "art", "magazines", "digest"]  # categories that are excluded and removed from data_names
76   file_filter = "^.*(\.(jpe?g|png|gif|bmp|zip|mp3|wav))|index\.html?$"  # use this to exclude by filenames
77   metadata_file_filter = "(x['metadata']['charratioB'] > 0.95) & " \
78                          "(x['metadata']['length'] > 300) & (x['metadata']['length'] < 30000)"
79   metadata_file_filter_declaration = "(x['metadata']['charratioB'] > 0.8) & " \
80                                      "(x['metadata']['length'] > 300) & (x['metadata']['length'] < 30000)"
81   data_dir = Path(project_path / "02_datasets/")
82   models_dir = Path(project_path / "03_workspace/models/")
83   analysis_dir = Path(project_path / "03_workspace/analysis/")
84   octis_dataset_dir = Path(project_path / "03_workspace/OCTIS/preprocessed_datasets/")
85   tmp_dir = Path(project_path / ".tmp/")
86
87   data_dir.mkdir(parents=True, exist_ok=True)
88   models_dir.mkdir(parents=True, exist_ok=True)
89   tmp_dir.mkdir(parents=True, exist_ok=True)
90
91   threads = []
92
93   if "download" not in skip_steps:
94       # prepare the threads for loading data
95       for names in helpers.chunker_list(data_names, number_of_threads):
96           threads.append(wrangler.loader(tmp_dir, data_dir, data_url, names))
97
```

```python
 98        # start all threads
 99        for thread in threads:
100            thread.start()
101
102        # wait for all threads to finish their work
103        running = True
104        while running:
105            running = False
106            for thread in threads:
107                if thread.is_alive():
108                    running = True
109                    print("waiting for threads to finish...")
110                    time.sleep(1)
111
112        print("data downloaded successfully")
113
114   if "cleaning" not in skip_steps:
115        # prepare the threads for cleaning up stuff
116        threads = []
117        dataset = {}
118
119        # prepare the threads
120        for names in helpers.chunker_list(data_names, number_of_threads):
121            threads.append(wrangler.cleaner(data_dir, names, file_filter))
122
123        # add declaration
124        threads.append(wrangler.cleaner(data_dir, ["declaration"], file_filter))
125
126        # start all threads
127        for thread in threads:
128            thread.start()
129
130        # wait for all threads to finish their work
131        running = True
132        while running:
133            running = False
134            for thread in threads:
135                if thread.is_alive():
136                    running = True
137                    print("waiting for threads to finish...")
138                    time.sleep(1)
139                else:
140                    if thread.data:
141                        dataset.update(thread.data)
142                        for d in thread.data:
143                            helpers.save_object(thread.data[d], tmp_dir.joinpath(str(d + ".pkl")))
144
145        helpers.save_object(dataset, tmp_dir.joinpath("dataset_full.pkl"))
146
147        # write metadata to csv file
148        print("Write dataset to csv")
149        f = open(tmp_dir.joinpath("dataset.csv"), "w+")
150        f.write("category,name,path,length,length_raw,avgcolumnsize,charratioA,charratioB,year,eyear,lyear,type\r\n")
151        for d in dataset:
152            for item in dataset[d]:
153                f.write("\"" + d + "\",\"" + str(item["metadata"]["name"]) + "\",\"" +
154                        str(item["metadata"]["path"]) + "\"," +
155                        str(item["metadata"]["length"]) + "," +
156                        str(item["metadata"]["length_raw"]) + "," +
157                        str(item["metadata"]["avgcolumnsize"]) + "," +
158                        str(item["metadata"]["charratioA"]) + "," +
159                        str(item["metadata"]["charratioB"]) + ",\"" +
160                        str(item["metadata"]["year"]) + "\"," +
161                        str(item["metadata"]["eyear"]) + "," +
162                        str(item["metadata"]["lyear"]) + ",\"" +
163                        str(item["metadata"]["type"]) + "\"\r\n")
164        f.close()
165
166        shutil.copyfile(tmp_dir.joinpath("dataset.csv"),
167                        models_dir.joinpath("dataset.csv")) # copy dataset.csv to models
168
169        print("data cleaned successfully")
170   elif "metadata-filtering" not in skip_steps:
171        # load dataset_full.pkl because it was not generate during runtime
172        dataset = helpers.load_object(tmp_dir.joinpath("dataset_full.pkl"))
173        print("data loaded from dataset_full.pkl")
174
175   if "metadata-filtering" not in skip_steps:
176        for key in data_names_exclude:
177            if key in dataset:
178                del dataset[key]
179
180        for key in dataset:
181            d_tmp = []
182            for x in dataset[key]:
183                if key == "declaration" and eval(metadata_file_filter_declaration):
184                    d_tmp.append(x)
185                elif eval(metadata_file_filter):
186                    d_tmp.append(x)
187            dataset[key] = d_tmp
188
189        helpers.save_object(dataset, tmp_dir.joinpath("dataset_filtered.pkl"))
190        shutil.copyfile(tmp_dir.joinpath("dataset_filtered.pkl"),
191                        models_dir.joinpath("dataset_filtered.pkl")) # copy dataset_filtered.pkl to models
192   elif "modeling" not in skip_steps:
193        # load dataset_filtered.pkl because it was not generate during runtime
194        dataset = helpers.load_object(tmp_dir.joinpath("dataset_filtered.pkl"))
195        print("data loaded from dataset_filtered.pkl")
196
197   if "modeling" not in skip_steps:
```

```python
198        t0 = time.time()
199        # calculate size of dataset
200        dataset_size = 0
201        for key in dataset:
202            dataset_size = dataset_size + len(dataset[key])
203
204        print("start building corpus")
205        BaseManager.register('PoolCorpus', nlp_pool.PoolCorpus)
206
207        if __name__ == '__main__':
208            with BaseManager() as manager:
209                corp = manager.PoolCorpus()
210                corp.set_totalFilesTarget(dataset_size)
211                with Pool(processes=number_of_threads) as pool:
212                    for key in dataset:
213                        pool.map(corp.add, ((d["content"], d["metadata"]) for d in dataset[key]))
214                corpus = corp.get()
215                print("corpus loaded")
216                corpus.save(tmp_dir.joinpath("corpus.bin.gz"))
217        print("end building corpus")
218        print("Time elapsed: ", time.time() - t0, "s")  # CPU seconds elapsed (floating point)
219
220        shutil.copyfile(tmp_dir.joinpath("corpus.bin.gz"),
221                        models_dir.joinpath("corpus.bin.gz")) # copy dataset_filtered.pkl to models
222
223    # load corpus.bin.gz always because freq cannot handle vanilla corpus
224    corpus = textacy.Corpus.load("en_core_web_sm", tmp_dir.joinpath("corpus.bin.gz"))
225    print("data loaded from corpus.bin.gz")
226
227    if "analysis_freq" not in skip_steps:
228        print("start wordcount")
229        # get lowercased and filtered corpus vocabulary (R3.3.1)
230        vocab = corpus.word_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
231        vocab_doc = corpus.word_doc_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
232
233        # sort vocabulary by descending frequency
234        vocab_sorted = sorted(vocab.items(), key=lambda x: x[1], reverse=True)
235        vocab_sorted_doc = sorted(vocab_doc.items(), key=lambda x: x[1], reverse=True)
236
237        # write to file, one word and its frequency per line
238        with open(tmp_dir.joinpath('vocab_frq.txt'), 'w') as f:
239            for word, frq in vocab_sorted:
240                line = f"{word}\t{frq}\n"
241                f.write(line)
242        with open(tmp_dir.joinpath('vocab_frq_doc.txt'), 'w') as f:
243            for word, frq in vocab_sorted_doc:
244                line = f"{word}\t{frq}\n"
245                f.write(line)
246
247        shutil.copyfile(tmp_dir.joinpath("vocab_frq.txt"),
248                        analysis_dir.joinpath("vocab_frq.txt")) # copy dataset_filtered.pkl to analysis
249        shutil.copyfile(tmp_dir.joinpath("vocab_frq_doc.txt"),
250                        analysis_dir.joinpath("vocab_frq_doc.txt"))  # copy dataset_filtered.pkl to analysis
251
252        en = textacy.load_spacy_lang("en_core_web_sm")
253        stats = []
254        for cat in data_names + ["declaration"]:
255            ctmp = corpus.get(lambda doc: doc._.meta["category"] == cat)
256            dtmp = list(ctmp)
257            ndocs = len(dtmp)
258            nvocab = 0
259            nlength = 0
260            for d in dtmp:
261                nvocab = nvocab + d.vocab.length
262                nlength = nlength + d._.meta["length"]
263            stats.append({"category": cat, "ndocs":  ndocs, "nvocab": nvocab, "nlength": nlength})
264            if ndocs > 0:
265                tmpcorpus = textacy.corpus.Corpus(en, data=dtmp)
266                tmpvocab = tmpcorpus.word_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
267                tmpvocab_doc = tmpcorpus.word_doc_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
268                tmpvocab_sorted = sorted(tmpvocab.items(), key=lambda x: x[1], reverse=True)
269                tmpvocab_sorted_doc = sorted(tmpvocab_doc.items(), key=lambda x: x[1], reverse=True)
270                with open(tmp_dir.joinpath('vocab_frq_' + cat + '.txt'), 'w') as f:
271                    for word, frq in tmpvocab_sorted:
272                        line = f"{word}\t{frq}\n"
273                        f.write(line)
274                with open(tmp_dir.joinpath('vocab_frq_doc_' + cat + '.txt'), 'w') as f:
275                    for word, frq in tmpvocab_sorted_doc:
276                        line = f"{word}\t{frq}\n"
277                        f.write(line)
278                shutil.copyfile(tmp_dir.joinpath('vocab_frq_' + cat + '.txt'),
279                                analysis_dir.joinpath('vocab_frq_' + cat + '.txt'))  # copy dataset_filtered.pkl to analysis
280                shutil.copyfile(tmp_dir.joinpath('vocab_frq_doc_' + cat + '.txt'),
281                                analysis_dir.joinpath('vocab_frq_doc_' + cat + '.txt'))  # copy dataset_filtered.pkl to analysis
282        with open(tmp_dir.joinpath('stats.csv'), 'w') as csv_file:
283            writer = csv.DictWriter(csv_file, stats[0].keys())
284            writer.writeheader()
285            writer.writerows(stats)
286        shutil.copyfile(tmp_dir.joinpath('stats.csv'),
287                        analysis_dir.joinpath('stats.csv'))  # copy stats.csv to analysis
288
289        print("end wordcount")
290
291    if "analysis_advance_preparation" not in skip_steps:
292        print("start advance preparation")
293        # merge metadata and actual content for each document in the corpus
294        # ugly, verbose syntax to merge two dictionaries
295        data = [{**doc._.meta, **{'text': doc.text}} for doc in corpus]
296
297        # create panda dataframe
```

```python
298        df = pd.DataFrame(data)
299
300        df_sub = df[(df['text'].str.len() > 10)]
301
302        # make new column containing all relevant metadata (showing in plot later on)
303        df_sub['descripton'] = df_sub[['name', 'year', 'charratioB', 'avgcolumnsize']].astype(str).agg(', '.join, axis=1)
304
305        helpers.save_object(df, tmp_dir.joinpath("df.pkl"))
306        helpers.save_object(df_sub, tmp_dir.joinpath("df_sub.pkl"))
307        print("end advance preparation")
308
309        shutil.copyfile(tmp_dir.joinpath("df.pkl"),
310                        models_dir.joinpath("df.pkl")) # copy df.pkl to models
311        shutil.copyfile(tmp_dir.joinpath("df_sub.pkl"),
312                        models_dir.joinpath("df_sub.pkl")) # copy df_sub.pkl to models
313    else:
314        # load df.pkl and df_sub.pkl because it was not generate during runtime
315        df = helpers.load_object(tmp_dir.joinpath("df.pkl"))
316        df_sub = helpers.load_object(tmp_dir.joinpath("df_sub.pkl"))
317        print("data loaded from df.pkl and df_sub.pkl")
318
319    if "analysis_scattertext" not in skip_steps:
320        print("start scattertext")
321        censor_tags = set(['CARD'])  # tags to ignore in corpus, e.g. numbers
322
323        en = textacy.load_spacy_lang("en_core_web_sm")
324        # stop words to ignore in corpus
325        en_stopwords = spacy.lang.en.stop_words.STOP_WORDS  # default stop words
326        custom_stopwords = set(['[', ']', '%'])
327        en_stopwords = en_stopwords.union(custom_stopwords)  # extend with custom stop words
328
329        # create corpus from dataframe
330        # lowercased terms, no stopwords, no numbers
331        # use lemmas for English only, German quality is too bad
332        corpus_speeches = st.CorpusFromPandas(df_sub,  # dataset
333                                              category_col='type',  # index differences by ...
334                                              text_col='text',
335                                              nlp=en,  # EN model
336                                              feats_from_spacy_doc=st.FeatsFromSpacyDoc(tag_types_to_censor=censor_tags,
337                                                                                        use_lemmas=True),
338                                              ).build().get_stoplisted_unigram_corpus(en_stopwords)
339
340        # produce visualization (interactive html)
341        html = st.produce_scattertext_explorer(corpus_speeches,
342                                               category='declaration',  # set attribute to divide corpus into two parts
343                                               category_name='declaration',
344                                               not_category_name='textfiles',
345                                               metadata=df_sub['descripton'],
346                                               width_in_pixels=1000,
347                                               minimum_term_frequency=5,  # drop terms occurring less than 5 times
348                                               save_svg_button=True,
349                                               )
350
351        # write visualization to html file
352        fname = tmp_dir.joinpath("viz_declaration_textfiles.html")
353        open(fname, 'wb').write(html.encode('utf-8'))
354        print("end scattertext")
355        shutil.copyfile(tmp_dir.joinpath("viz_declaration_textfiles.html"),
356                        analysis_dir.joinpath("viz_declaration_textfiles.html")) # copy viz_declaration_textfiles.html to analysis
357
358    if "analysis_year" not in skip_steps:
359        print("start year")
360
361        df_sub = df[(df['text'].str.len() > 10)]
362
363        # make new column containing all relevant metadata (showing in plot later on)
364        df_sub['descripton'] = df_sub[['name', 'year', 'charratioB', 'avgcolumnsize']].astype(str).agg(', '.join, axis=1)
365
366        dtmp = df_sub.groupby('eyear').agg({'text': "count"}).reset_index().rename(columns={'text': 'count'})
367        dtmp = dtmp.rename(columns={"eyear": "year"})
368        dtmp.insert(2, "type", "eyear")
369        docs_per_year = dtmp
370
371        dtmp = df_sub.groupby('lyear').agg({'text': "count"}).reset_index().rename(columns={'text': 'count'})
372        dtmp = dtmp.rename(columns={"lyear": "year"})
373        dtmp.insert(2, "type", "lyear")
374        docs_per_year = docs_per_year.append(dtmp, ignore_index=True)
375
376        # manual year was only available in top100 analysis
377        # dtmp = pd.read_csv('top100_years.txt', delimiter=",").groupby('myear').agg({'text': "count"}).reset_index().rename(
378        #     columns={'text': 'count'})
379        # dtmp = dtmp.rename(columns={"myear": "year"})
380        # dtmp.insert(2, "type", "myear")
381        # docs_per_year = docs_per_year.append(dtmp, ignore_index=True)
382
383        docs_per_year = docs_per_year[docs_per_year["year"] != "NA"]
384        docs_per_year['year'] = pd.to_numeric(docs_per_year['year'])
385
386        p = (ggplot(docs_per_year, aes('year', 'count', color='type', group='type'))
387             + geom_point(alpha=0.5, stroke=0)
388             + geom_line()
389             + theme_classic()
390             + labs(x="Year",
391                    y="absolute number",
392                    color="Legend")
393             + theme(axis_text_x=element_text(angle=90, hjust=1))
394             + scale_x_continuous(limits=(1960, 1999))
395             )
396
397        ggsave(plot=p, filename="docs_per_year", path=tmp_dir)
```

```
398
399        shutil.copyfile(tmp_dir.joinpath("docs_per_year.png"),
400                        analysis_dir.joinpath("docs_per_year.png"))  # copy docs_per_year.png to analysis
401
402        try:
403            dummy = pd.DataFrame(years.from_1960_to_1999)
404
405            e = dummy.append(docs_per_year[(docs_per_year['type'] == "eyear") &
406                             (docs_per_year['year'] >= 1960)][["year", "count"]])\
407                .groupby('year').agg({'count': "sum"})
408            l = dummy.append(docs_per_year[(docs_per_year['type'] == "lyear") &
409                             (docs_per_year['year'] >= 1960)][["year", "count"]])\
410                .groupby('year').agg({'count': "sum"})
411
412            print("r_{eyear mit lyear} = ", np.corrcoef(e["count"], l["count"])[0, 1])
413        except Exception as e:
414            print("something went wrong while calculating eyear / lyear variaty.")
415            pass
416
417        print("end year")
418
419    if "analysis_octis" not in skip_steps:
420        print("start octis")
421
422        df_sub.to_csv(tmp_dir.joinpath("corpus.txt"), "\t", columns = ["text"])
423
424        # Initialize preprocessing
425        preprocessor = Preprocessing(vocabulary=None, max_features=None,
426                                     remove_punctuation=True, punctuation=string.punctuation,
427                                     lemmatize=True, stopword_list='english',
428                                     min_chars=1, min_words_docs=0)
429
430        # preprocess
431        octis_dataset = preprocessor.preprocess_dataset(documents_path=tmp_dir.joinpath("corpus.txt"))
432
433        # save the preprocessed dataset
434        octis_dataset.save(str(tmp_dir.joinpath('octis_dataset')))
435
436        if os.path.exists(octis_dataset_dir.joinpath('octis_dataset')):
437            shutil.rmtree(octis_dataset_dir.joinpath('octis_dataset'))
438        if os.path.exists(models_dir.joinpath('octis_dataset')):
439            shutil.rmtree(models_dir.joinpath('octis_dataset'))
440
441        shutil.copytree(tmp_dir.joinpath('octis_dataset'),
442                        octis_dataset_dir.joinpath('octis_dataset'))  # copy octis_dataset to datasets
443        shutil.copytree(tmp_dir.joinpath('octis_dataset'),
444                        models_dir.joinpath('octis_dataset'))  # copy octis_dataset to models
445
446        # octis will be processed using dashboard
447        # model = LDA(num_topics=5)  # Create model
448        # model_output = model.train_model(octis_dataset)  # Train the model
449
450        # metric = TopicDiversity(topk=10)  # Initialize metric
451        # topic_diversity_score = metric.score(model_output)  # Compute score of the metric
452
453        # print("topic diversity score:", topic_diversity_score)
454
455        # wordcloud will be generated outside of python
456        # wocl = WordCloud(mode="RGBA", background_color="white").generate(" ".join(model_output["topics"][0]))
457        # wocl2 = WordCloud(mode="RGBA", background_color="white", relative_scaling=1, scale=10, max_words=9999,
458        #                   min_font_size=1, max_font_size=18, collocations=False).generate(" ".join(model_output["topics"][0]))
459        # image = wocl.to_image()
460        # image.save(tmp_dir.joinpath("wordcloud.png"))
461        # image2 = wocl2.to_image()
462        # image2.save(tmp_dir.joinpath("wordcloud2.png"))
463
464        print("end octis")
465
466    if "analysis_entities" not in skip_steps:
467        entities = []
468
469        for doc in corpus.docs:
470            for ent in textacy.extract.entities(doc):
471                try:
472                    entities += [{"text": ent.text, "label": ent.label_, "explain": spacy.explain(ent.label_)}]
473                except:
474                    print("Problem with:", doc._.meta["name"])
475
476        # export corpus as csv
477        f_csv = tmp_dir.joinpath('entities.csv')
478        textacy.io.csv.write_csv(entities, f_csv, fieldnames=entities[0].keys())
479
480        shutil.copyfile(tmp_dir.joinpath('entities.csv'),
481                        analysis_dir.joinpath('entities.csv'))  # copy entities.csv to analysis
482
483        df_entities = pd.DataFrame(entities, columns=['text', 'label', 'explain'])
484        df_entities_count = df_entities.groupby('text').agg({'label': "count"}).rename(
485            columns={'label': 'count'}).sort_values(by=['count'], ascending=False).reset_index()
486
487        # write to file, one word and its frequency per line
488        fname = tmp_dir.joinpath('entities_frq.csv')
489        with open(fname, 'w') as f:
490            for i, d in df_entities_count.iterrows():
491                line = d["text"] + "," + str(d["count"]) + "\n"
492                f.write(line)
493
494        shutil.copyfile(tmp_dir.joinpath('entities_frq.csv'),
495                        analysis_dir.joinpath('entities_frq.csv'))  # copy entities_frq.csv to analysis
496
497        print("entities: \n", df_entities_count[:25])
498
499    print("everything done.")
```

```python
1   #!/usr/bin/env python3
2   # -*- coding: utf-8 -*-
3
4   ########################################################
5   # title:   datawrangler
6   # author: Josias Bruderer
7   # date:    27.08.2021
8   # desc:    this module takes care of all tasks that are
9   #          related to juggling files and datasets.
10  ########################################################
11
12  import os
13  import sys
14  from threading import Thread
15  from pathlib import Path
16  import requests
17  import zipfile
18  import codecs
19  import re
20
21
22  def averageLen(lst, excludeEmpty=True):
23      if excludeEmpty:
24          lengths = [len(i) for i in lst if i != "" and i != " "]
25      else:
26          lengths = [len(i) for i in lst]
27      return 0 if len(lengths) == 0 else round((float(sum(lengths)) / len(lengths)), 2)
28
29
30  def daterange(lst, t="r"):
31      ltmp = []
32      ltmp2 = []
33      if len(lst) > 0:
34          for l in lst:
35              ltmp += list(filter(None, l))
36          for l in ltmp:
37              if len(l) == 2:
38                  ltmp2 += ["19" + l]
39              else:
40                  ltmp2 += [l]
41          if len(ltmp2) > 2:
42              if t == "e":
43                  return str(min(ltmp2))
44              elif t == "l":
45                  return str(max(ltmp2))
46              else:
47                  return str(min(ltmp2) + "-" + max(ltmp2))
48          else:
49              return str(ltmp2[0])
50      else:
51          return "NA"
52
53  class cleaner(Thread):
54
55      def __init__(self, data_dir, data_names, file_filter):
56          Thread.__init__(self)
57          self.data_dir = data_dir
58          self.data_names = data_names
59          self.file_filter = file_filter
60          self.data = {}
61
62      def run(self):
63          for data_name in self.data_names:
64              self.data[data_name] = self.get_texts(Path(self.data_dir, data_name), data_name)
65
66      def get_texts(self, dir_texts, data_name):
67          """
68          Sequentially stream all documents from a given folder,
69          including metadata.
70          """
71          data = []
72
73          # iterate over all documents
74          for fname in dir_texts.glob('**/*'):  # ** = all subdirectories
75              if Path(fname).is_file():
76                  print("processing in " + str(dir_texts.stem) + " file: " + str(fname))
77
78                  if re.match(self.file_filter, fname.name):
79                      print("skip in " + str(dir_texts.stem) + " because of file_filter: " + str(fname))
80                      continue
81                  # Read file content and replace encoding erros
82                  content_raw = codecs.open(fname, 'r', encoding='utf-8', errors='replace').read()
83
84                  # join lines as there are hard line-breaks
85                  content = content_raw.replace('\r\n', ' ')
86                  content = content.replace('\r', ' ')
87                  content = content.replace('\n', ' ')
```

```
 88                      content = content.replace('\t', ' ')
 89                      content = content.replace('\x1a', ' ')
 90                      content = re.sub('[^A-z0-9\ \.\'\,\!]', ' ', content)
 91                      content = re.sub('[\\\\\^\[\]]', ' ', content)
 92
 93                      # add more metadata here if needed
 94                      # charratio: 0 = no character is "text", 1 = every character is "text"
 95                      if len(content_raw) == 0:
 96                          charratioA = 0
 97                          charratioB = 0
 98                      else:
 99                          charratioA = round(1 / len(content_raw) * len(re.findall("[A-z]", content_raw)), 2)
100                          charratioB = round(1 / len(content_raw) * len(re.findall("[A-z\ \.\"\,\!]", content_raw)), 2)
101
102                      typ = "textfile"
103                      if fname.name == "declarationbarlow1996.txt" or data_name == "declaration":
104                          typ = "declaration"
105
106                      rxdate = re.compile(
107                          'copyright.{0,3}(19[6-9][0-9])|updated.{0,3}[0-1]?[0-9]?-[0-3]?[0-9]?-([6-9][0-9])|'
108                          'Date\:.*([6-9][0-9]).*,|(?:jan(?:uary)?|feb(?:ruary)?|mar(?:ch)?|apr(?:il)?|may|'
109                          'june|july|aug(?:ust)?|sept(?:ember)?|oct(?:ober)?|nov(?:ember)?|dec(?:ember)?)'
110                          '.{0,8}(1?9?[6-9][0-9])|[0-1]?[0-9]?\/[0-3]?[0-9]?\/([6-9][0-9])|[0-1]?[0-9]?-[0-3]?[0-9]?'
111                          '-([6-9][0-9])|[^-](19[6-9][0-9])')
112
113                      matches = rxdate.findall(content, re.IGNORECASE)
114
115                      metadata = {'name': fname.name,
116                                  'path': str(fname),
117                                  'length_raw': len(content_raw),
118                                  'length': len(content),
119                                  'avgcolumnsize': averageLen(content_raw.splitlines()),
120                                  'charratioA': charratioA,
121                                  'charratioB': charratioB,
122                                  'year': daterange(matches),
123                                  'eyear': daterange(matches, "e"),
124                                  'lyear': daterange(matches, "l"),
125                                  'type': typ,
126                                  'category': data_name,
127                                  }
128
129                      # return documents one after another (sequentially)
130                      data.append({"content": content, "metadata": metadata})
131          return data
132
133  class loader(Thread):
134
135      def __init__(self, tmp_dir, data_dir, data_url, data_names):
136          Thread.__init__(self)
137          self.tmp_dir = tmp_dir
138          self.data_dir = data_dir
139          self.data_url = data_url
140          self.data_names = data_names
141          self.init()
142
143      def init(self):
144          try:
145              # create tmp directory if not existing yet
146              if not os.path.exists(self.tmp_dir.is_dir()):
147                  os.mkdir(self.tmp_dir)
148
149              # create data directory if not existing yet
150              if not os.path.exists(self.data_dir):
151                  os.mkdir(self.data_dir)
152          except:
153              print("Unexpected error: ", sys.exc_info()[0])
154              raise
155
156      def run(self):
157          for data_name in self.data_names:
158              url = self.data_url.replace("[name]", data_name)
159              zipdir = Path(self.tmp_dir, str(data_name + ".zip"))
160              self.download_zip(url, zipdir)
161              self.extract_zip(zipdir)
162
163      def download_zip(self, url, zipdir):
164          try:
165              if not zipdir.is_file():
166                  print("Downloading: " + url)
167                  r = requests.get(url, allow_redirects=True)
168                  open(zipdir, 'wb').write(r.content)
169              else:
170                  print("Skip downloading, file already downloaded: " + url)
171          except:
172              print("Unexpected error: ", sys.exc_info()[0])
173              raise
174
175      def extract_zip(self, zipdir):
176          try:
177              if not Path(self.data_dir / zipdir.stem).is_dir():
```

```
178                    print("Extracting: " + str(zipdir))
179                    with zipfile.ZipFile(zipdir, 'r') as zip_ref:
180                        zip_ref.extractall(self.data_dir)
181                else:
182                    print("Skip extracting, file already extracted: " + str(zipdir))
183            except:
184                print("Unexpected error: ", sys.exc_info()[0])
185                raise
186
```

```python
 1   #!/usr/bin/env python3
 2   # -*- coding: utf-8 -*-
 3
 4   ########################################################
 5   # title:   helpers
 6   # author: Josias Bruderer
 7   # date:    25.08.2021
 8   # desc:    this module provides some useful functions
 9   ########################################################
10
11   import pickle
12
13
14   def save_object(obj, filename):
15       with open(filename, 'wb') as outp:  # Overwrites any existing file.
16           pickle.dump(obj, outp, pickle.HIGHEST_PROTOCOL)
17
18   def load_object(filename):
19       with open(filename, 'rb') as inp:
20           return pickle.load(inp)
21
22   def chunker_list(seq, size):
23       return (seq[i::size] for i in range(size))
24
```

```python
1    #!/usr/bin/env python3
2    # -*- coding: utf-8 -*-
3
4    ######################################################
5    # title:  helpers
6    # author: Josias Bruderer
7    # date:   28.08.2021
8    # desc:   this module provides nlp functions
9    ######################################################
10
11   import textacy
12   import spacy
13   # run: ./.envs/bin/python -m spacy download en_core_web_sm
14
15
16   class PoolCorpus(object):
17
18       def __init__(self):
19           model = spacy.load('en_core_web_sm', disable=["parser"])
20           model.max_length = 10000000  # enable utilization of ~ 100GB RAM
21           self.corpus = textacy.corpus.Corpus(lang=model)
22           self.totalFilesTarget = 1
23           self.processedFiles = 0
24
25       def add(self, data):
26           self.corpus.add(data)
27           self.processedFiles = self.processedFiles + 1
28           print("Processed ", self.processedFiles, " of ", self.totalFilesTarget, "files: ",
29                 round(100/self.totalFilesTarget*self.processedFiles, 4), "%")
30
31       def get(self):
32           return self.corpus
33
34       def save(self, path):
35           self.corpus.save(path)
36
37       def set_totalFilesTarget(self, n):
38           self.totalFilesTarget = n
39
40   """
41   texts = {
42           'key1': 'First text 1.',
43           'key2': 'Second text 2.',
44           'key3': 'Third text 3.',
45           'key4': 'Fourth text 4.',
46       }
47
48   BaseManager.register('PoolCorpus', PoolCorpus)
49
50   if __name__ == '__main__':
51       with BaseManager() as manager:
52           corpus = manager.PoolCorpus()
53
54           with Pool(processes=2) as pool:
55               pool.map(corpus.add, ((v, {'key': k}) for k, v in texts.items()))
56
57           print(corpus.get())
58   """
```

```python
#!/usr/bin/env python3
# -*- coding: utf-8 -*-

###########################################################
# title:  helpers
# author: Josias Bruderer
# date:   29.08.2021
# desc:   this module provides years variable
###########################################################

from_1960_to_1999 = [{"year": 1960, "count": 0},
                     {"year": 1961, "count": 0},
                     {"year": 1962, "count": 0},
                     {"year": 1963, "count": 0},
                     {"year": 1964, "count": 0},
                     {"year": 1965, "count": 0},
                     {"year": 1966, "count": 0},
                     {"year": 1967, "count": 0},
                     {"year": 1968, "count": 0},
                     {"year": 1969, "count": 0},
                     {"year": 1970, "count": 0},
                     {"year": 1971, "count": 0},
                     {"year": 1972, "count": 0},
                     {"year": 1973, "count": 0},
                     {"year": 1974, "count": 0},
                     {"year": 1975, "count": 0},
                     {"year": 1976, "count": 0},
                     {"year": 1977, "count": 0},
                     {"year": 1978, "count": 0},
                     {"year": 1979, "count": 0},
                     {"year": 1980, "count": 0},
                     {"year": 1981, "count": 0},
                     {"year": 1982, "count": 0},
                     {"year": 1983, "count": 0},
                     {"year": 1984, "count": 0},
                     {"year": 1985, "count": 0},
                     {"year": 1986, "count": 0},
                     {"year": 1987, "count": 0},
                     {"year": 1988, "count": 0},
                     {"year": 1989, "count": 0},
                     {"year": 1990, "count": 0},
                     {"year": 1991, "count": 0},
                     {"year": 1992, "count": 0},
                     {"year": 1993, "count": 0},
                     {"year": 1994, "count": 0},
                     {"year": 1995, "count": 0},
                     {"year": 1996, "count": 0},
                     {"year": 1997, "count": 0},
                     {"year": 1998, "count": 0},
                     {"year": 1999, "count": 0}]
```