

main.py

```
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 analysis 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", "piracy",
75                       "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) & \" \" \\\n"
78                       "(x['metadata']['length'] > 300) & (x['metadata']['length'] < 30000)"
79 data_dir = Path(project_path / "02_datasets/")
80 models_dir = Path(project_path / "03_workspace/models/")
81 analysis_dir = Path(project_path / "03_workspace/analysis/")
82 octis_dataset_dir = Path(project_path / "03_workspace/OCTIS/preprocessed_datasets/")
83 tmp_dir = Path(project_path / ".tmp/")
84
85 data_dir.mkdir(parents=True, exist_ok=True)
86 models_dir.mkdir(parents=True, exist_ok=True)
87 tmp_dir.mkdir(parents=True, exist_ok=True)
88
89 threads = []
90
91 if "download" not in skip_steps:
92     # prepare the threads for loading data
93     for names in helpers.chunker_list(data_names, number_of_threads):
```

```

94         threads.append(wrangler.loader(tmp_dir, data_dir, data_url, names))
95
96     # start all threads
97     for thread in threads:
98         thread.start()
99
100    # wait for all threads to finish their work
101    running = True
102    while running:
103        running = False
104        for thread in threads:
105            if thread.is_alive():
106                running = True
107                print("waiting for threads to finish...")
108                time.sleep(1)
109
110    print("data downloaded successfully")
111
112    if "cleaning" not in skip_steps:
113        # prepare the threads for cleaning up stuff
114        threads = []
115        dataset = {}
116
117        # prepare the threads
118        for names in helpers.chunker_list(data_names, number_of_threads):
119            threads.append(wrangler.cleaner(data_dir, names, file_filter))
120
121        # add declaration
122        threads.append(wrangler.cleaner(data_dir, ["declaration"], file_filter))
123
124        # start all threads
125        for thread in threads:
126            thread.start()
127
128        # wait for all threads to finish their work
129        running = True
130        while running:
131            running = False
132            for thread in threads:
133                if thread.is_alive():
134                    running = True
135                    print("waiting for threads to finish...")
136                    time.sleep(1)
137                else:
138                    if thread.data:
139                        dataset.update(thread.data)
140                        for d in thread.data:
141                            helpers.save_object(thread.data[d], tmp_dir.joinpath(str(d + ".pkl")))
142
143            helpers.save_object(dataset, tmp_dir.joinpath("dataset_full.pkl"))
144
145        # write metadata to csv file
146        print("Write dataset to csv")
147        f = open(tmp_dir.joinpath("dataset.csv"), "w+")
148        f.write("category,name,path,length,length_raw,avgcolumnsize,charratioA,charratioB,year,eyear,lyear,type\r\n")
149        for d in dataset:
150            for item in dataset[d]:
151                f.write("\\" + d + "\\", \"\" + str(item["metadata"]["name"]) + "\",\"\" +
152                    str(item["metadata"]["path"]) + "\",\"\" +
153                    str(item["metadata"]["length"]) + "\",\"\" +
154                    str(item["metadata"]["length_raw"]) + "\",\"\" +
155                    str(item["metadata"]["avgcolumnsize"]) + "\",\"\" +
156                    str(item["metadata"]["charratioA"]) + "\",\"\" +
157                    str(item["metadata"]["charratioB"]) + "\",\"\" +
158                    str(item["metadata"]["year"]) + "\",\"\" +
159                    str(item["metadata"]["eyear"]) + "\",\"\" +
160                    str(item["metadata"]["lyear"]) + "\",\"\" +
161                    str(item["metadata"]["type"]) + "\"\r\n")
162            f.close()
163
164        shutil.copyfile(tmp_dir.joinpath("dataset.csv"),
165                        models_dir.joinpath("dataset.csv")) # copy dataset.csv to models
166
167        print("data cleaned successfully")
168    elif "metadata-filtering" not in skip_steps:
169        # load dataset_full.pkl because it was not generate during runtime
170        dataset = helpers.load_object(tmp_dir.joinpath("dataset_full.pkl"))
171        print("data loaded from dataset_full.pkl")
172
173    if "metadata-filtering" not in skip_steps:
174        for key in data_names_exclude:
175            if key in dataset:
176                del dataset[key]
177
178        for key in dataset:
179            d_tmp = []
180            for x in dataset[key]:
181                if eval(metadata_file_filter):
182                    d_tmp.append(x)
183            dataset[key] = d_tmp
184
185        helpers.save_object(dataset, tmp_dir.joinpath("dataset_filtered.pkl"))
186        shutil.copyfile(tmp_dir.joinpath("dataset_filtered.pkl"),
187                        models_dir.joinpath("dataset_filtered.pkl")) # copy dataset_filtered.pkl to models
188    elif "modeling" not in skip_steps:
189        # load dataset_filtered.pkl because it was not generate during runtime

```

```

190 dataset = helpers.load_object(tmp_dir.joinpath("dataset_filtered.pkl"))
191 print("data loaded from dataset_filtered.pkl")
192
193 if "modeling" not in skip_steps:
194     t0 = time.time()
195     # calculate size of dataset
196     dataset_size = 0
197     for key in dataset:
198         dataset_size = dataset_size + len(dataset[key])
199
200     print("start building corpus")
201     BaseManager.register('PoolCorpus', nlp_pool.PoolCorpus)
202
203     if __name__ == '__main__':
204         with BaseManager() as manager:
205             corp = manager.PoolCorpus()
206             corp.set_totalFilesTarget(dataset_size)
207             with Pool(processes=number_of_threads) as pool:
208                 for key in dataset:
209                     pool.map(corp.add, ((d["content"], d["metadata"]) for d in dataset[key]))
210             corpus = corp.get()
211             print("corpus loaded")
212             corpus.save(tmp_dir.joinpath("corpus.bin.gz"))
213         print("end building corpus")
214         print("Time elapsed: ", time.time() - t0, "s") # CPU seconds elapsed (floating point)
215
216         shutil.copyfile(tmp_dir.joinpath("corpus.bin.gz"),
217                         models_dir.joinpath("corpus.bin.gz")) # copy dataset_filtered.pkl to models
218
219 # load corpus.bin.gz always because freq cannot handle vanilla corpus
220 corpus = textacy.Corpus.load("en_core_web_sm", tmp_dir.joinpath("corpus.bin.gz"))
221 print("data loaded from corpus.bin.gz")
222
223 if "analysis_freq" not in skip_steps:
224     print("start wordcount")
225     # get lowercased and filtered corpus vocabulary (R3.3.1)
226     vocab = corpus.word_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
227     vocab_doc = corpus.word_doc_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
228
229     # sort vocabulary by descending frequency
230     vocab_sorted = sorted(vocab.items(), key=lambda x: x[1], reverse=True)
231     vocab_sorted_doc = sorted(vocab_doc.items(), key=lambda x: x[1], reverse=True)
232
233     # write to file, one word and its frequency per line
234     with open(tmp_dir.joinpath('vocab_frq.txt'), 'w') as f:
235         for word, frq in vocab_sorted:
236             line = f"{word}\t{frq}\n"
237             f.write(line)
238     with open(tmp_dir.joinpath('vocab_frq_doc.txt'), 'w') as f:
239         for word, frq in vocab_sorted_doc:
240             line = f"{word}\t{frq}\n"
241             f.write(line)
242
243     shutil.copyfile(tmp_dir.joinpath("vocab_frq.txt"),
244                     analysis_dir.joinpath("vocab_frq.txt")) # copy dataset_filtered.pkl to analysis
245     shutil.copyfile(tmp_dir.joinpath("vocab_frq_doc.txt"),
246                     analysis_dir.joinpath("vocab_frq_doc.txt")) # copy dataset_filtered.pkl to analysis
247
248     en = textacy.load_spacy_lang("en_core_web_sm")
249     stats = []
250     for cat in data_names + ["declaration"]:
251         ctmp = corpus.get(lambda doc: doc._meta["category"] == cat)
252         dtmp = list(ctmp)
253         ndocs = len(dtmp)
254         nvocab = 0
255         nlength = 0
256         for d in dtmp:
257             nvocab = nvocab + d.vocab.length
258             nlength = nlength + d._meta["length"]
259         stats.append({"category": cat, "ndocs": ndocs, "nvocab": nvocab, "nlength": nlength})
260     if ndocs > 0:
261         tmpcorpus = textacy.corpus.Corpus(en, data=dtmp)
262         tmpvocab = tmpcorpus.word_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
263         tmpvocab_doc = tmpcorpus.word_doc_counts(by='lemma_', filter_stops=True, filter_punct=True, filter_nums=True)
264         tmpvocab_sorted = sorted(tmpvocab.items(), key=lambda x: x[1], reverse=True)
265         tmpvocab_sorted_doc = sorted(tmpvocab_doc.items(), key=lambda x: x[1], reverse=True)
266         with open(tmp_dir.joinpath('vocab_frq_' + cat + '.txt'), 'w') as f:
267             for word, frq in tmpvocab_sorted:
268                 line = f"{word}\t{frq}\n"
269                 f.write(line)
270         with open(tmp_dir.joinpath('vocab_frq_doc_' + cat + '.txt'), 'w') as f:
271             for word, frq in tmpvocab_sorted_doc:
272                 line = f"{word}\t{frq}\n"
273                 f.write(line)
274         shutil.copyfile(tmp_dir.joinpath('vocab_frq_' + cat + '.txt'),
275                         analysis_dir.joinpath('vocab_frq_' + cat + '.txt')) # copy dataset_filtered.pkl to analysis
276         shutil.copyfile(tmp_dir.joinpath('vocab_frq_doc_' + cat + '.txt'),
277                         analysis_dir.joinpath('vocab_frq_doc_' + cat + '.txt')) # dataset_filtered.pkl to analysis
278     with open(tmp_dir.joinpath('stats.csv'), 'w') as csv_file:
279         writer = csv.DictWriter(csv_file, stats[0].keys())
280         writer.writeheader()
281         writer.writerows(stats)
282
283     print("end wordcount")
284
285 if "analysis_advance_preparation" not in skip_steps:

```

```

286 print("start advance preparation")
287 # merge metadata and actual content for each document in the corpus
288 # ugly, verbose syntax to merge two dictionaries
289 data = [{**doc._meta, **{'text': doc.text}} for doc in corpus]
290
291 # create panda dataframe
292 df = pd.DataFrame(data)
293
294 df_sub = df[(df['text'].str.len() > 10)]
295
296 # make new column containing all relevant metadata (showing in plot later on)
297 df_sub['descripton'] = df_sub[['name', 'year', 'charratioB', 'avgcolumnsize']].astype(str).agg(' '.join, axis=1)
298
299 helpers.save_object(df, tmp_dir.joinpath("df.pkl"))
300 helpers.save_object(df_sub, tmp_dir.joinpath("df_sub.pkl"))
301 print("end advance preparation")
302
303 shutil.copyfile(tmp_dir.joinpath("df.pkl"),
304                 models_dir.joinpath("df.pkl")) # copy df.pkl to models
305 shutil.copyfile(tmp_dir.joinpath("df_sub.pkl"),
306                 models_dir.joinpath("df_sub.pkl")) # copy df_sub.pkl to models
307 else:
308 # load df.pkl and df_sub.pkl because it was not generate during runtime
309 df = helpers.load_object(tmp_dir.joinpath("df.pkl"))
310 df_sub = helpers.load_object(tmp_dir.joinpath("df_sub.pkl"))
311 print("data loaded from df.pkl and df_sub.pkl")
312
313 if "analysis_scattertext" not in skip_steps:
314 print("start scattertext")
315 censor_tags = set(['CARD']) # tags to ignore in corpus, e.g. numbers
316
317 en = textacy.load_spacy_lang("en_core_web_sm")
318 # stop words to ignore in corpus
319 en_stopwords = spacy.lang.en.stop_words.STOP_WORDS # default stop words
320 custom_stopwords = set(['', ' ', '%'])
321 en_stopwords = en_stopwords.union(custom_stopwords) # extend with custom stop words
322
323 # create corpus from dataframe
324 # lowercased terms, no stopwords, no numbers
325 # use lemmas for English only, German quality is too bad
326 corpus_speeches = st.CorpusFromPandas(df_sub, # dataset
327                                     category_col='type', # index differences by ...
328                                     text_col='text',
329                                     nlp=en, # EN model
330                                     feats_from_spacy_doc=st.FeatsFromSpacyDoc(tag_types_to_censor=censor_tags,
331                                     use_lemmas=True),
332                                     ).build().get_stoplisted_unigram_corpus(en_stopwords)
333
334 # produce visualization (interactive html)
335 html = st.produce_scattertext_explorer(corpus_speeches,
336                                     category='declaration', # set attribute to divide corpus into two parts
337                                     category_name='declaration',
338                                     not_category_name='textfiles',
339                                     metadata=df_sub['descripton'],
340                                     width_in_pixels=1000,
341                                     minimum_term_frequency=5, # drop terms occurring less than 5 times
342                                     save_svg_button=True,
343                                     )
344
345 # write visualization to html file
346 fname = tmp_dir.joinpath("viz_declaration_textfiles.html")
347 open(fname, 'wb').write(html.encode('utf-8'))
348 print("end scattertext")
349 shutil.copyfile(tmp_dir.joinpath("viz_declaration_textfiles.html"),
350                 analysis_dir.joinpath("viz_declaration_textfiles.html")) # viz_declaration_textfiles.html >analysis
351
352 if "analysis_year" not in skip_steps:
353 print("start year")
354
355 df_sub = df[(df['text'].str.len() > 10)]
356
357 # make new column containing all relevant metadata (showing in plot later on)
358 df_sub['descripton'] = df_sub[['name', 'year', 'charratioB', 'avgcolumnsize']].astype(str).agg(' '.join, axis=1)
359
360 dtmp = df_sub.groupby('eyear').agg({'text': "count"}).reset_index().rename(columns={'text': 'count'})
361 dtmp = dtmp.rename(columns={'eyear': "year"})
362 dtmp.insert(2, "type", "eyear")
363 docs_per_year = dtmp
364
365 dtmp = df_sub.groupby('lyear').agg({'text': "count"}).reset_index().rename(columns={'text': 'count'})
366 dtmp = dtmp.rename(columns={'lyear': "year"})
367 dtmp.insert(2, "type", "lyear")
368 docs_per_year = docs_per_year.append(dtmp, ignore_index=True)
369
370 # manual year was only available in top100 analysis
371 # dtmp = pd.read_csv('top100_years.txt', delimiter=",").groupby('myear').agg({'text': "count"}).reset_index()
372 # .rename(columns={'text': 'count'})
373 # dtmp = dtmp.rename(columns={'myear': "year"})
374 # dtmp.insert(2, "type", "myear")
375 # docs_per_year = docs_per_year.append(dtmp, ignore_index=True)
376
377 docs_per_year = docs_per_year[docs_per_year["year"] != "NA"]
378 docs_per_year['year'] = pd.to_numeric(docs_per_year['year'])
379
380 p = (ggplot(docs_per_year, aes('year', 'count', color='type', group='type'))
381      + geom_point(alpha=0.5, stroke=0))

```

```

382         + geom_line()
383         + theme_classic()
384         + labs(x="Year",
385               y="absolute number",
386               color="Legend")
387         + theme(axis_text_x=element_text(angle=90, hjust=1))
388         + scale_x_continuous(limits=(1960, 1999))
389     )
390
391     ggsave(plot=p, filename="docs_per_year", path=tmp_dir)
392
393     shutil.copyfile(tmp_dir.joinpath("docs_per_year.png"),
394                     analysis_dir.joinpath("docs_per_year.png")) # copy docs_per_year.png to analysis
395
396     try:
397         dummy = pd.DataFrame(years.from_1960_to_1999)
398
399         e = dummy.append(docs_per_year[(docs_per_year['type'] == "eyear") &
400                                         (docs_per_year['year'] >= 1960)][["year", "count"]])\
401             .groupby('year').agg({'count': "sum"})
402         l = dummy.append(docs_per_year[(docs_per_year['type'] == "lyear") &
403                                         (docs_per_year['year'] >= 1960)][["year", "count"]])\
404             .groupby('year').agg({'count': "sum"})
405
406         print("r_{eyear mit lyear} = ", np.corrcoef(e["count"], l["count"])[0, 1])
407     except Exception as e:
408         print("something went wrong while calculating eyear / lyear variaty.")
409         pass
410
411     print("end year")
412
413     if "analysis_octis" not in skip_steps:
414         print("start octis")
415
416         df_sub.to_csv(tmp_dir.joinpath("corpus.txt"), "\t", columns = ["text"])
417
418         # Initialize preprocessing
419         preprocessor = Preprocessing(vocabulary=None, max_features=None,
420                                     remove_punctuation=True, punctuation=string.punctuation,
421                                     lemmatize=True, stopwords_list='english',
422                                     min_chars=1, min_words_docs=0)
423
424         # preprocess
425         octis_dataset = preprocessor.preprocess_dataset(documents_path=tmp_dir.joinpath("corpus.txt"))
426
427         # save the preprocessed dataset
428         octis_dataset.save(str(tmp_dir.joinpath('octis_dataset')))
429
430         if os.path.exists(octis_dataset_dir.joinpath('octis_dataset')):
431             shutil.rmtree(octis_dataset_dir.joinpath('octis_dataset'))
432         if os.path.exists(models_dir.joinpath('octis_dataset')):
433             shutil.rmtree(models_dir.joinpath('octis_dataset'))
434
435         shutil.copytree(tmp_dir.joinpath('octis_dataset'),
436                         octis_dataset_dir.joinpath('octis_dataset')) # copy octis_dataset to datasets
437         shutil.copytree(tmp_dir.joinpath('octis_dataset'),
438                         models_dir.joinpath('octis_dataset')) # copy octis_dataset to models
439
440         # octis will be processed using dashboard
441         # model = LDA(num_topics=5) # Create model
442         # model_output = model.train_model(octis_dataset) # Train the model
443
444         # metric = TopicDiversity(topk=10) # Initialize metric
445         # topic_diversity_score = metric.score(model_output) # Compute score of the metric
446
447         # print("topic diversity score:", topic_diversity_score)
448
449         # wordcloud will be generated outside of python
450         # wocl = WordCloud(mode="RGBA", background_color="white").generate(" ".join(model_output["topics"][0]))
451         # wocl2 = WordCloud(mode="RGBA", background_color="white", relative_scaling=1, scale=10, max_words=9999,
452                             # min_font_size=1, max_font_size=18, collocations=False).generate(" "
453                                                 # .join(model_output["topics"][0]))
454         # image = wocl.to_image()
455         # image.save(tmp_dir.joinpath("wordcloud.png"))
456         # image2 = wocl2.to_image()
457         # image2.save(tmp_dir.joinpath("wordcloud2.png"))
458
459         print("end octis")
460
461     if "analysis_entities" not in skip_steps:
462         entities = []
463
464         for doc in corpus.docs:
465             for ent in textacy.extract.entities(doc):
466                 try:
467                     entities += [{"text": ent.text, "label": ent.label_, "explain": spacy.explain(ent.label_)}]
468                 except:
469                     print("Problem with:", doc._meta["name"])
470
471         # export corpus as csv
472         f_csv = tmp_dir.joinpath('entities.csv')
473         textacy.io.csv.write_csv(entities, f_csv, fieldnames=entities[0].keys())
474
475         shutil.copyfile(tmp_dir.joinpath('entities.csv'),
476                         analysis_dir.joinpath('entities.csv')) # copy entities.csv to analysis
477

```

```

478 df_entities = pd.DataFrame(entities, columns=['text', 'label', 'explain'])
479 df_entities_count = df_entities.groupby('text').agg({'label': "count"}).rename(
480     columns={'label': 'count'}).sort_values(by=['count'], ascending=False).reset_index()
481
482 # write to file, one word and its frequency per line
483 fname = tmp_dir.joinpath('entities_frq.csv')
484 with open(fname, 'w') as f:
485     for i, d in df_entities_count.iterrows():
486         line = d["text"] + "," + str(d["count"]) + "\n"
487         f.write(line)
488
489 shutil.copyfile(tmp_dir.joinpath('entities_frq.csv'),
490                 analysis_dir.joinpath('entities_frq.csv')) # copy entities_frq.csv to analysis
491
492 print("entities: \n", df_entities_count[:25])
493
494 print("everything done.")

```

modules/helpers.py

```
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
```

modules/nlp_pool.py

```

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.Corpora(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      """

```



```

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":
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|june|'
109             '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
```

```
1  #!/usr/bin/env python3
2  # -*- coding: utf-8 -*-
3
4  #####
5  # title:  helpers
6  # author: Josias Bruderer
7  # date:   29.08.2021
8  # desc:   this module provides years variable
9  #####
10
11 from_1960_to_1999 = [{"year": 1960, "count": 0},
12                      {"year": 1961, "count": 0},
13                      {"year": 1962, "count": 0},
14                      {"year": 1963, "count": 0},
15                      {"year": 1964, "count": 0},
16                      {"year": 1965, "count": 0},
17                      {"year": 1966, "count": 0},
18                      {"year": 1967, "count": 0},
19                      {"year": 1968, "count": 0},
20                      {"year": 1969, "count": 0},
21                      {"year": 1970, "count": 0},
22                      {"year": 1971, "count": 0},
23                      {"year": 1972, "count": 0},
24                      {"year": 1973, "count": 0},
25                      {"year": 1974, "count": 0},
26                      {"year": 1975, "count": 0},
27                      {"year": 1976, "count": 0},
28                      {"year": 1977, "count": 0},
29                      {"year": 1978, "count": 0},
30                      {"year": 1979, "count": 0},
31                      {"year": 1980, "count": 0},
32                      {"year": 1981, "count": 0},
33                      {"year": 1982, "count": 0},
34                      {"year": 1983, "count": 0},
35                      {"year": 1984, "count": 0},
36                      {"year": 1985, "count": 0},
37                      {"year": 1986, "count": 0},
38                      {"year": 1987, "count": 0},
39                      {"year": 1988, "count": 0},
40                      {"year": 1989, "count": 0},
41                      {"year": 1990, "count": 0},
42                      {"year": 1991, "count": 0},
43                      {"year": 1992, "count": 0},
44                      {"year": 1993, "count": 0},
45                      {"year": 1994, "count": 0},
46                      {"year": 1995, "count": 0},
47                      {"year": 1996, "count": 0},
48                      {"year": 1997, "count": 0},
49                      {"year": 1998, "count": 0},
50                      {"year": 1999, "count": 0}]
51
```