

```

1 # -*- coding: utf-8 -*-
2 import subprocess
3 import re
4 import io
5
6 PDFFILES="pdffiles.txt"
7 PDFPREFIX="pdf/"
8 TXTPREFIX="txt/"
9 OUTFILE="matches.xls"
10
11 YEARS=("2008", "2009", "2010", "2011", "2012", "2013", "2014", "2015", "2016", "2017", "2018")
12
13 PATTERNS={"ECOLOGY_TERM"      :{"ecology"      : r"\b[Ee]colog[A-Za-z]+\b"},
14         "ECOLOGY_MODE"      :{"trophic1"      : r"\b[A-Za-z-]*troph[A-Za-z]*\b",
15         "trophic2"      : r"\b[A-Za-z-]+trophic\b",
16         "saprobe"      : r"\bsaprob[A-Za-z]*\b",
17         "mutualis"     : r"\bmutualis[A-Za-z]+\b",
18         "parasite"         : r"\bparasit[A-Za-z]+\b",
19         "commensal"    : r"\bcommensa[A-Za-z]+\b",
20         "mycorrhiza"   : r"\b[A-Za-z-]*mycorrhiz[A-Za-z]+\b",
21         "phyte"        : r"\b[A-Za-z-]+phyte[A-Za-z]*\b",
22         "phytic"       : r"\b[A-Za-z-]+phytic[A-Za-z]*\b",
23         "predator"        : r"\b[A-Za-z-]*predator[A-Za-z]*\b",
24         "predacious"  : r"\b[A-Za-z-]*predaceous[A-Za-z]*\b",
25         "biont"        : r"\b[A-Za-z-]+biont[A-Za-z]*\b",
26         "pathogen"        : r"\b[A-Za-z-]*pathog[A-Za-z]+\b",
27         "lithic"       : r"\b[A-Za-z-]+lithic[A-Za-z]*\b",
28         "endo"         : r"\b[A-Za-z-]*endo[A-Za-z]+\b"},
29         "ECOLOGY_ASSOC"    :{"substrate"     : r"\b[A-Za-z-]*substrates*\b",
30         "habitat"         : r"\b[A-Za-z-]*habitats*\b",
31         "host"           : r"\b[A-Za-z-]*hosts*\b",
32         "partner"        : r"\b[A-Za-z-]*partners*\b"},
33         "COUNTRY"         :{"country"       : r"\bcountre[A-Za-z-]+\b"},
34         "LOCALITY"        :{"local"        : r"\blocalit[A-Za-z-]+\b"},
35         "GIS/GPS"         :{"GIS"          : r"\bgis\b",
36         "GPS"            : r"\bgps\b",
37         "coord"       : r"\b[A-Za-z-]+ordinates*\b"},
38         "DISTRIBUTION"   :{"geograph"   : r"\bgeograph[A-Za-z]+\b",
39         "dist"           : r"\bdistributed[A-Za-z]+\b"},
40         "ALTITUDE"        :{"alt"          : r"\baltitud[A-Za-z]+\b",
41         "ele"         : r"\belevati[A-Za-z]+\b"},
42         "CLIMATE_ZONE"   :{"boreal"       : r"\b[A-Za-z-]*boreal\b",
43         "temper"         : r"\b[A-Za-z-]*temperate\b",
44         "arctic"         : r"\b[A-Za-z-]*arctic[A-Za-z]*\b",
45         "tropic"         : r"\b[A-Za-z-]*tropic[A-Za-z]*\b",
46         "africa"     : r"\b[A-Za-z-]*africa[A-Za-z]*\b",
47         "europa"     : r"\b[A-Za-z-]*europa[A-Za-z]*\b",
48         "america"    : r"\b[A-Za-z-]*america[A-Za-z]*\b",
49         "asia"       : r"\b[A-Za-z-]*asia[A-Za-z]*\b",
50         "austral"        : r"\baustral[A-Za-z]*\b"},
51         "CLIMATE"         :{"climate"      : r"\bclimat[A-Za-z]+\b"},
52         "MOLECULAR_MARKER":{"DNA"         : r"\b[A-Za-z-]*dna\b",
53         "seq"           : r"\bsequenc[A-Za-z]+\b",
54         "PCR"           : r"\b[A-Za-z-]*pcr\b"},
55         "MOLECULAR_DB"   :{"genbank"     : r"\bgenbank\b",
56         "INSD"          : r"\binsd\b",
57         "INSDC"         : r"\binsdc\b",
58         "ENA"           : r"\bena\b",
59         "EMBL"          : r"\bembl\b",
60         "DDBJ"          : r"\bddbj\b"},
61         "MOLECULAR_AVAIL":{"treebase"    : r"\btreebase\b",
62         "dryad"         : r"\bdryad\b"},
63         "MOLECULAR_BLAST":{"blast"       : r"\b[tnpx]*blast[tnpx]*\b"},
64         "PHYLUM"         :{"mycota"      : r"\b[A-Za-z]+mycota\b"},
65         "ORDER"          :{"ales"       : r"\b[A-Za-z]{4,50}ales\b"},
66         "FAMILY"         :{"aceae"      : r"\b[A-Za-z]+aceae\b"},
67         "BIODIVERSITY"   :{"biodiv"     : r"\b\b"},
68         "THREATEND"      :{"threat"     : r"\bthreatened\b",
69         "end"           : r"\bendangered\b"},

```

```

70         "COLLECTION"      : {"herb"       : r"\bherbari[A-Za-z]+\b",
71                             "fung"       : r"\bfungari[A-Za-z]+\b",
72                             "museum"     : r"\bmuseum[A-Za-z]*\b",
73                             "collection": r"\bculture collection[s]*\b"},
74         "SUPPLEMENTAL"   : {"supp1"      : r"\bsupplementa[A-Za-z]+\b"},
75         "INDEX_FUNGORUM" : {"if"        : r"\bindex fungorum\b",
76                             "iflink"     : r"indexfungorum"},
77         "MYCOBANK"       : {"mb"         : r"\bmycobank\b"},
78         "FUNGAL_NAMES"   : {"fn"         : r"\bfungal names\b"},
79         "SOCIAL_IMPL"    : {"agri"       : r"\bagricultur[A-Za-z]+\b",
80                             "for"        : r"\bforestry\b",
81                             "biotech"    : r"\bbiotechnolog[A-Za-z]+\b",
82                             "med"        : r"\bmedic[A-Za-z]+\b",
83                             "clinic"     : r"\bclinic[A-Za-z]+\b",
84                             "health"     : r"\bhealth *care\b",
85                             "ind"        : r"\bindustr[A-Za-z]+\b",
86                             "aqua"       : r"\baquacultur[A-Za-z]+\b",
87                             "horti"      : r"\bhorticultur[A-Za-z]+\b"}}
88
89     EXCLUDEDPATTERNS={"ECOLOGY_MODE": {"trophic1" : r"electrophor",
90                                         "endo"    : r"ependorf"}}
91
92
93     def readPDFList(PDFFILE):
94
95         rawdata=open(PDFFILE).readlines()
96
97         pdffiles=[]
98
99         for line in rawdata:
100             pdffiles.append(line.strip())
101
102         return pdffiles
103
104
105     def convertPDFtoTXT(pdffiles):
106
107         txtfiles=[]
108
109         for pdffile in pdffiles:
110
111             print "Converting %s" % pdffile
112
113             txtfile=pdffile+".txt"
114             txtfile_old=pdffile+"-old.txt"
115
116             process=subprocess.Popen(["pdftotext", PDFPREFIX+pdffile, TXTPREFIX+txtfile])
117             stdout, stderr = process.communicate()
118
119             process=subprocess.Popen(["pdftotext", "-raw", PDFPREFIX+pdffile,
120                                     TXTPREFIX+txtfile_old])
121             stdout, stderr = process.communicate()
122
123             txtfiles.append(txtfile)
124
125         return txtfiles
126
127     def processArticles(txtlist):
128
129         all_matches={}
130
131         for txtfile in txtlist:
132             print "Parsing %s" % txtfile
133             all_matches[txtfile]=parseArticle(TXTPREFIX+txtfile)
134
135         return all_matches
136
137     def parseArticle(txtfile):

```

```

138 matches={}
139 references_found=0
140 abstract_found=0
141 acknowledgements_found=0
142 funding_found=0
143
144 print "Parsing %s" % txtfile
145
146 indata=open(txtfile).readlines()
147
148 cntr_start=0
149
150 # Search for year
151
152 # This code search for a year associated with several key phrases
153 # Exceptions is made for Mycologia since those PDFs often contains an additional page
154 # where the
155 # available online date can differ significantly from the published date. Also, in some
156 # Mycologia
157 # papers, the years appears on the same line and for some the year appears on the line
158 # below.
159 # The loop continutes until a year has been found. If no year is found, the value will
160 # be -1.
161
162 myco_year_on_next_line=0
163 is_mycologia=0
164 year=-1
165
166 # Check first line. Some papers have spaces between the numbers.
167 tmp=indata[0].replace(" ", "")
168 for yr in YEARS:
169     if yr in tmp:
170         year=yr
171         break
172
173 # Continue if no year was found.
174 if year==-1:
175     for row in indata:
176         row=row.strip().rstrip().lower()
177
178         if year!=-1:
179             break
180
181         if myco_year_on_next_line==1:
182             for yr in YEARS:
183                 if yr in row:
184                     year=yr
185                     break
186
187         if row[:21]=="to cite this article":
188             for yr in YEARS:
189                 if yr in row:
190                     year=yr
191                     break
192
193         if "accepted:" in row:
194             tmp=row.split("accepted:")[1].split(" ")[1].strip()
195             if tmp in YEARS:
196                 year=tmp
197                 break
198
199         if row[:16]=="available online" and is_mycologia==0:
200             for yr in YEARS:
201                 if yr in row:
202                     year=yr
203                     break
204
205         if row[:9]=="mycologia":
206             is_mycologia=1

```

```

203         for yr in YEARS:
204             if yr in row:
205                 year=yr
206                 break
207         if year==-1:
208             myco_year_on_next_line=1
209
210     if year==-1:
211         print "Warning: Year not found."
212     else:
213         print "Found year %s" % year
214
215     matches["year"]=year
216
217     # Search for abstract
218     for row in indata:
219
220         row=row.strip().lower()
221         row=row.replace(" ", "")
222
223         if row=="abstract" or row=="a b s t r a c t" or row[:8]=="abstract" or row[:15]=="a
224             b s t r a c t":
225             abstract_found=1
226             print "Found abstract at row %d" % cntr_start
227             break
228
229         cntr_start+=1
230
231     if abstract_found==0:
232         print "Warning: No abstract found."
233
234     cntr=0
235     tot_row=len(indata)
236     for row in indata:
237
238         if cntr>=cntr_start and (("the mycological society of america" not in row.lower())
239             or ("copyright" not in row.lower())):
240             row=row.strip().lower()
241             for section in PATTERNS:
242                 for pattern in PATTERNS[section]:
243                     foundmatch=0
244                     if re.search(PATTERNS[section][pattern], row)!=None:
245                         if section in EXCLUDEDPATTERNS and pattern in
246                             EXCLUDEDPATTERNS[section]:
247                             if re.search(EXCLUDEDPATTERNS[section][pattern], row)==None:
248                                 matches[section]=1
249                                 foundmatch=1
250                     else:
251                         matches[section]=1
252                         foundmatch=1
253                     if foundmatch==1:
254                         print "Found %s (%s, %s) in %s." % (pattern, section,
255                             PATTERNS[section][pattern], row)
256
257     #
258     print 1.0*cntr/tot_row, row
259     if (((1.0*cntr)/tot_row)>0.50 and (row=="references" or row=="r e f e r e n c e
260         s" or row=="literature cited" or row=="l i t e r a t u r e   c i t e d" or
261         row=="reference" or row=="list of references")):
262         references_found=1
263         break
264
265     if (((1.0*cntr)/tot_row)>0.50 and (row=="acknowledgements" or row=="a c k n o w
266         l e d g e m e n t s" or row=="acknowledgments" or row=="a c k n o w l e d g m e
267         n t s" or row[:16]=="acknowledgements" or row[:15]=="acknowledgments")):
268         acknowledgements_found=1
269         break

```

```

264         if (((1.0*cntr)/tot_row)>0.50 and (row=="funding" or row=="f u n d i n g")):
265             funding_found=1
266             break
267
268
269         cntr+=1
270
271     if acknowledgements_found==1:
272         print "Exited at acknowledgements (row %d of %d)" % (cntr, tot_row)
273     else:
274         print "No acknowledgements found"
275
276     if references_found==0:
277         print "Warning: No references found!"
278     else:
279         print "Exited at reference (row %d of %d)" % (cntr, tot_row)
280
281     if funding_found==0:
282         print "Warning: No funding found!"
283     else:
284         print "Exited at funding (row %d or %d)" % (cntr, tot_row)
285
286     print "Done!\n\n"
287     return matches
288
289 def writeResults(matches):
290
291     colnames=sorted(PATTERNS.keys())
292     # print colnames
293
294     of=open(OUTFILE, "w")
295
296     of.write("Papers\tJournal\tYear\t%s\n" % "\t".join(colnames))
297
298     for paper in matches:
299         curline=[]
300         for colname in colnames:
301             if colname in matches[paper]:
302                 curline.append("Y")
303             else:
304                 curline.append("N")
305
306         tmp=paper.split("/")
307         journal=tmp[-2]
308         article=tmp[-1]
309
310         of.write("%s\t%s\t%s\t%s\n" % (article, journal,
311             matches[paper]["year"], "\t".join(curline)))
312
313     of.close()
314
315     return
316
317
318
319
320 pdflist=readPDFList(PDFFILES)
321 txtlist=convertPDFtoTXT(pdflist)
322 #txtlist=["warpeace.txt"]
323 matches=processArticles(txtlist)
324 writeResults(matches)
325

```