

chap2023.py

```
1 # cos類似度と極性分類の確信度から類似文書を表示する。知識情報演習3に独自の工夫を追加。
2 # クエリに対して文書のスコアを出すプログラム(0 <= SCORE <= 1)
3
4
5 import math
6 import pandas as pd
7 import re
8 from janome.tokenizer import Tokenizer
9 from asari.api import Sonar
10 import os
11 from IPython.terminal.prompts import Token
12 import stopwords
13
14
15 sonar = Sonar()
16 t = Tokenizer()
17
18 query = '日本 ワールドカップ'
19 #query = 'ナゲツ ファイナル'
20 #query = 'NBA トレード'
21 #query = '渡邊雄太 NBA'
22 #query = 'ウェンバンヤマ ドラフト'
23 argv = query.split()
24 argc = len(argv)
25
26 # クエリ表示
27 print(argv)
28 # ストップワードリストを取得
29 stoplist = stopwords.stopwords()
30
31 # cos類似度を計算 -----
32 def cos(argv):
33     index_file = "index/index.txt"
34
35     idf_scores = {}
36     tfidf_scores = {}
37     query_tf = {}
38     query_tfidf = {}
39     query_words = {}
40     ranking_docs = {}
41
42     #tfidf_scoresからデータフレームを作成
43     f = open(index_file, 'r')
44
45     for line in f:
46         line = line.rstrip()
47         split_line = line.split('\t')
48         word = split_line[0]
49         doc = split_line[1]
50         idf = float(split_line[2])
51         tfidf = float(split_line[3])
52
53         if word not in tfidf_scores:
54             tfidf_scores[word] = {}
55             tfidf_scores[word][doc] = tfidf
56         else:
```

```
57         tfidf_scores[word][doc] = tfidf
58
59     if word not in idf_scores:
60         idf_scores[word] = {}
61         idf_scores[word] = idf
62     else:
63         idf_scores[word] = idf
64
65     tfidf_table = pd.DataFrame(tfidf_scores)
66     tfidf_table = tfidf_table.fillna(0)
67
68
69 #クエリに対する処理
70 argv = '吾輩は猫である'
71 query_file = 'query'
72
73 pattern = re.compile(r'^[ --]$')
74
75
76 #tokens = t.tokenize(argv[1])
77 tokens = []
78 for i in range(argc):
79     list = t.tokenize(argv[i])
80     tokens.append(list)
81 flat_t = [item for sublist in tokens for item in sublist]
82
83 flat_token = []
84 for token in flat_t:
85     tmp = token.surface
86     judge = pattern.match(tmp)
87     flat_token.append(tmp)
88
89     if judge:
90         continue
91     if tmp in stoplist:
92         continue
93
94     if tmp in query_words:
95         query_words[tmp] += 1
96     else:
97         query_words[tmp] = 1
98
99
100
101 #クエリtfidf_scoresからクエリのデータフレームを作る
102 for index_word in idf_scores:
103     query_tf[index_word] = {}
104     query_tfidf[index_word] = {}
105     query_tf[index_word][query_file] = 0
106     query_tfidf[index_word][query_file] = 0
107
108     for query_word in query_words:
109         if query_word in query_tf:
110             query_tf[query_word][query_file] += 1
111
112         for query_word in query_words:
113             if query_word in query_tf:
114                 query_tfidf[query_word][query_file] = float(query_tf[query_word][query_file]) * float(idf_scores[query_word][query_file])
```

```

116     query_table = pd.DataFrame(query_tfidf)
117
118     # クエリが索引後辞書の中でどの文書名に属すのかを確認し、対象文書を同定
119     for query_word in query_words:
120         if query_word in tfidf_scores:
121             for doc in tfidf_scores[query_word]:
122                 ranking_docs[doc] = 1
123
124     # cos類似度計算
125     for doc in ranking_docs:
126         numerator = 0
127         doc_vec = tfidf_table.loc[doc]
128         query_vec = query_table.loc[query_file]
129         for i in range(len(query_vec.values)):
130             i_value = query_vec.values[i] * doc_vec.values[i]
131             numerator = numerator + i_value
132
133         denominator = 0
134         query_value = 0
135         doc_value = 0
136         for i in range(len(query_vec.values)):
137             query_value += query_vec.values[i]**2
138             doc_value += doc_vec.values[i]**2
139
140         denominator = math.sqrt(query_value) * math.sqrt(doc_value)
141         cosine = numerator / denominator
142
143         ranking_docs[doc] = cosine
144
145     return ranking_docs.items() #sorted(ranking_docs.items(), key=lambda x:x[1],
reverse=True)
146
147
148
149 # 文書のポジティブ/ネガティブを変数scaledで0から1で表示 -----
150 def text_emotions():
151     data = []
152     for filename in os.listdir('text'):
153         doc_box = []
154         f = open('text/' + filename, 'r')
155         for row in f:
156             row = row.rstrip()
157             doc_box.append(row)
158         result_doc = '\n'.join(doc_box)
159         res = sonar.ping(text=result_doc)
160         classes = res['classes']
161         positive_confidence = next((item['confidence'] for item in classes if
item['class_name'] == 'positive'), 0)
162         negative_confidence = next((item['confidence'] for item in classes if
item['class_name'] == 'negative'), 0)
163
164         scaled = (positive_confidence - negative_confidence + 1) / 2
165         data.append(scaled)
166     return data
167
168
169 # クエリのポジティブ/ネガティブを変数scaledで-1から1で表示 -----
170 def query_emotions(argv):
171     query = ' '.join(argv)
172     res = sonar.ping(text=query)
173     classes = res['classes']

```

```
174     positive_confidence = next((item['confidence'] for item in classes if
175         item['class_name'] == 'positive'), 0)
176     negative_confidence = next((item['confidence'] for item in classes if
177         item['class_name'] == 'negative'), 0)
178     scaled = (positive_confidence - negative_confidence + 1) / 2
179     return scaled
180
181 # ポジネガの検索式を定義 -----
182 # クエリと文書のポジネガが違かったら、距離を調整(0.5倍)
183 def dis_emotions(qemo, temo):
184     if qemo >= 0.5:
185         if temo >= 0.5:
186             dis = abs(qemo - temo)
187         else:
188             dis = abs(qemo - temo) * 0.5
189     elif qemo < 0.5:
190         if temo <= 0.5:
191             dis = abs(qemo - temo)
192         else:
193             dis = abs(qemo - temo) * 0.5
194     scaled = 1 - dis*2
195     return scaled
196
197
198
199 # 検索スコアのためのアルゴリズムを定義 -----
200 # score = cos類似度*1/3 + ポジネガ確信度*1/3
201 # 0 < score < 1 に調整
202 def score(cos, emo):
203     return cos*(2/3) + emo*(1/3)
204
205
206
207
208 # クエリと文書のポジネガの距離を計算 -----
209 te = text_emotions()
210 qe = query_emotions(argv)
211
212 emo_dis = []
213 counter = 1
214 for i in te:
215     tmp = dis_emotions(qe, i)
216     emo_dis.append(['doc{}'.format(counter), tmp])
217     counter += 1
218
219
220 # 表示 -----
221 score_pair = []
222 cosinfo = cos(argv)
223 coslist = list(cosinfo)
224 for i in range(len(cosinfo)):
225     doc_score = score(coslist[i][1], emo_dis[i][1])
226     score_pair.append(['doc{}'.format(i+1), doc_score])
227
228 for i in sorted(score_pair, key=lambda x: x[1], reverse=True):
229     print(i)
```