

```

1 from gensim.models.coherencemodel import CoherenceModel
2 import pandas as pd
3 import re
4 from gensim import corpora, models
5 from gensim.parsing.preprocessing import preprocess_string
6 import warnings
7 import matplotlib.pyplot as plt
8 import matplotlib.style as style
9 style.use('fivethirtyeight')
10 warnings.filterwarnings('ignore')
11
12 pd.set_option('display.max_columns', None)
13 # pd.set_option('display.width', None)
14 # pd.set_option('display.max_colwidth', None)
15 pd.set_option('display.unicode.ambiguous_as_wide', True)
16 pd.set_option('display.unicode.east_asian_width', True)
17 # 引入上次已經斷詞的小王子資料且依照章節進行區分document
18 # 更新部分：去除長度為一的字母和某些表示數量的詞彙（中文）以及某些可能較為沒意義的詞彙
19 # list = ('一', '二', '三', '四', '五', '六', '七', '八', '九', '十', '兩',
20 #         '幾', '小王子')
21 # '是', '只', '那', '在', '再', '說', '這', '第')開頭的字(除了章節名稱之後會用到)
22 # 由於程式碼過於冗長就不附上
23
24 # first section
25
26 # data = open('output.txt', 'r', encoding='utf8')
27 # document = ['']*27
28 # tempword = []
29 # countchapters = ["第一章 \n", "第二章 \n",
30 #                   "第三章 \n", "第四章 \n", "第五章 \n", "第六章 \n",
31 #                   "第七章 \n", "第八章 \n", "第九章 \n", "第十章 \n",
32 #                   "第十一章 \n", "第十二章 \n", "第十三章 \n", "第十四章 \n",
33 #                   "第十五章 \n", "第十六章 \n", "第十七章 \n", "第十八章 \n",
34 #                   "第十九章 \n", "第二十章 \n", "第二十一章 \n", "第二十二章 \n",
35 #                   "第二十三章 \n", "第二十四章 \n",
36 #                   "第二十五章 \n", "第二十六章 \n", "第二十七章 \n", "作者 \n"]
37
38 # for line in data.readlines():
39 #     tempword.append(line) # 每行加入tempwords
40 # print(len(tempword))
41
42 # chapter = 0
43 # for i in range(len(tempword)):
44 #     if chapter < 27:
45 #         if tempword[i] == countchapters[chapter]:
46 #             temp = i+2
47 #             while True:
48 #                 document[chapter] += tempword[temp]
49 #                 temp += 1
50 #                 if tempword[temp] == countchapters[chapter+1]:
51 #                     f = open("allchapter/chapter"+str(chapter+1)+".txt",
52 #                             "w+")
53 #                     document[chapter] = re.sub("\n", " ",
54 # document[chapter])
55 #                     str_list = document[chapter].split()
56 #                     new_str = ' '.join(str_list)
57 #                     f.write(new_str)
58 #                     i = temp
59 #                     break

```

```
56 #             chapter += 1
57
58 # second section
59
60 # df = pd.DataFrame(columns=['document', 'text'])
61 # for i in range(27):
62 #     temprow = pd.read_table("chapter"+str(i+1)+".txt", header=None)
63 #     df = df.append(pd.DataFrame(
64 #         {'document': i, 'text': [temprow.iloc[0][0]]}),
65 #         ignore_index=True)
66
67 # df.to_csv("outputcsv.csv", index=False)
68 # print(df.text)
69
70 # third section
71
72 df = pd.read_csv('outputcsv.csv')
73
74 train_data = []
75 for i in range(27):
76     train_data.append((df.iloc[i, 1]).split())
77
78 def create_lsa_model(documents, dictionary, number_of_topics):
79     print(f'Creating LDA Model with {number_of_topics} topics')
80     document_terms = [dictionary.doc2bow(doc) for doc in documents]
81     return models.LsiModel(document_terms,
82                             num_topics=number_of_topics,
83                             id2word=dictionary,)
84
85
86 def run_lsa_process(documents, number_of_topics=10):
87     dictionary = corpora.Dictionary(documents)
88     lsa_model = create_lsa_model(documents, dictionary,
89                                 number_of_topics)
90     return documents, dictionary, lsa_model
91
92
93 def calculate_coherence_score(documents, dictionary, model):
94     coherence_model = CoherenceModel(model=model,
95                                     texts=documents,
96                                     dictionary=dictionary,
97                                     coherence='c_v')
98     return coherence_model.get_coherence()
99
100
101 def get_coherence_values(start, stop):
102     for num_topics in range(start, stop):
103         print(f'\nCalculating coherence for {num_topics} topics')
104         documents, dictionary, model = run_lsa_process(train_data,
105                                                         number_of_topics=num_topics)
106         coherence = calculate_coherence_score(documents,
107                                             dictionary,
108                                             model)
109         yield coherence
110
111
112 if __name__ == '__main__':
113
```

```
114     min_topics, max_topics = 1, 28
115     coherence_scores = list(get_coherence_values(min_topics, max_topics))
116     print("The highests coherence score: "+str(max(coherence_scores))+ " " +
117           "The best number of topics:
118 "+str(coherence_scores.index(max(coherence_scores))+1))
119     x = [int(i) for i in range(min_topics, max_topics)]
120     documents, dictionary, model = run_lsa_process(
121         train_data, coherence_scores.index(max(coherence_scores))+1)
122
123     topics = pd.DataFrame(columns=['index', '主題內容'])
124
125     finaltopics = model.print_topics(
126         num_topics=coherence_scores.index(max(coherence_scores))+1,
127         num_words=10)
128     for topic in finaltopics:
129         topics = topics.append(pd.DataFrame(
130             {'index': [topic[0]+1], '主題內容': [topic[1]]}))
131     topics.to_csv("topics.csv", index=False)
132     plt.figure(figsize=(10, 8))
133     plt.ylim([0, 1])
134     plt.xticks(x, x)
135     plt.plot(x, coherence_scores, marker='o')
136     plt.xlabel('Number of topics')
137     plt.ylabel('Coherence Value')
138     plt.title('Coherence Scores by number of Topics')
139
140     plt.annotate("Best"+"\\n"+str(max(coherence_scores))+ "\\n"+str(coherence_sco
141 res.index(max(coherence_scores))+1),
142                 xy=(coherence_scores.index(max(coherence_scores))+1,
143 max(coherence_scores)), xycoords='data')
144     plt.savefig("outputplt.jpg")
145
```