

▼ B0928022 杜云驊

Colab-link: https://colab.research.google.com/drive/1q4XG_TBgLN2U8nvAfrMk-xdhjYz3LFrU?usp=sharing

```
1 import requests
2 import json
3 import jieba
4 import networkx as nx
5 import re
6 from bs4 import BeautifulSoup
7
8 # • 15064
9 MOVIE_TOTAL_NUM = 15064
10 LINK = "https://movies.yahoo.com.tw/movieinfo_main/"
11
12
13 class MovieInfo:
14     def __init__(self, doc_id=0, cname='unknown', ename='unknown', pagerank=0, label='unknown', intro='unknown', re
15         links=[]) -> None:
16         self.doc_id = doc_id
17         self.cname = cname
18         self.ename = ename
19         self.pagerank = pagerank
20         self.label = label
21         self.intro = intro
22         self.released_date = released_date
23         self.links = list(set([14653, 13733, 14652, 14690, 14848, 14301, 14558, 14822, 14850, 14557, 14208,
24
25     def generate_dict(self) -> dict:
26         dict = {}
27         dict['doc_id'] = self.doc_id
28         dict['cname'] = self.cname
29         dict['ename'] = self.ename
30         dict['pagerank'] = self.pagerank
31         dict['label'] = self.label
32         dict['intro'] = self.intro
33         dict['released_date'] = self.released_date
34         dict['links'] = self.links
35
36     return dict
37
38
39 class SearchEngine:
40     def __init__(self, inverted_index, path) -> None:
41         self.inverted_index = inverted_index
42         self.movie_data = []
43         self._get_movie_data(path)
44
45     def query(self, query) -> None:
46         seg = jieba.cut_for_search(query)
47         words = {word for word in seg}
48         match_results = set(self.inverted_index[words.pop()])
49         while words:
50             match_results.intersection_update(self.inverted_index[words.pop()])
51
52         match_results = self._sort_by_rankpage(match_results)
53
54         if match_results == set():
55             print("None")
56         self._print_query_result(match_results, query)
57
58     def _get_movie_data(self, path) -> None:
59         with open(path) as f:
60             self.movie_data = json.load(f)
61
62     def _print_query_result(self, id_list, query) -> None:
63         print(f"您的搜尋結果:")
64         print(f"共{len(id_list)}筆, 符合\"{query}\"")
65
66         regex = re.compile(rf'(.*){query}')
```

```

30         relate = 0
31         for movie in self.movie_data:
32             if regex.match(movie['cname']) or regex.match(movie['ename']) or regex.match(movie['intro']):
33                 relate += 1
34
35         count = 0
36         for id in id_list:
37             for movie in self.movie_data:
38                 if movie['doc_id'] == id:
39                     if regex.match(movie['cname']) or regex.match(movie['ename']) or regex.match(movie['
40                         count += 1
41                         print(f"\nid: {movie['pagerank']}"):")
42                         print(f"Chinese Name: {movie['cname']}")
43                         print(f"English Name: {movie['ename']}")
44                         print(f"Intro: {movie['intro']}")
45         print(f"\nPrecision = {count} / {len(id_list)} = {count / len(id_list) * 100}%")
46         print(f"Recall = {count} / {relate} = {count / relate * 100}%")
47
48     def _sort_by_rankpage(self, matches) -> list:
49         dict = {}
50         for id in matches:
51             for movie in self.movie_data:
52                 if movie['doc_id'] == id:
53                     dict[id] = movie['pagerank']
54                     break
55
56         return [key for key, pagerank in sorted(dict.items(), key=lambda x: x[1], reverse=True)]
57

```

```

1 response = requests.get("https://movies.yahoo.com.tw/chart.html")
2 soup = BeautifulSoup(response.text, "lxml")
3
4 resault = soup.find('div', 'rank_list table rankstyle1').find('dl', 'rank_list_box').find('h2').text.strip()
5 lt1 = [resault]
6 resault = soup.find('div', 'rank_list table rankstyle1').find_all('div', 'rank_txt')[:9]
7 for m in resault:
8     lt1.append(m.text.strip())
9
10 response = requests.get("https://movies.yahoo.com.tw/chart.html?cate=us")
11 soup = BeautifulSoup(response.text, "lxml")
12
13 resault = soup.find('div', 'rank_list table rankstyle1').find('dl', 'rank_list_box').find('h2').text.strip()
14 lt2 = [resault]
15 resault = soup.find('div', 'rank_list table rankstyle1').find_all('div', 'rank_txt')[:9]
16 for m in resault:
17     lt2.append(m.text.strip())
18
19 response = requests.get("https://movies.yahoo.com.tw/chart.html?cate=trailer")
20 soup = BeautifulSoup(response.text, "lxml")
21
22 resault = soup.find('div', 'rank_list table rankstyle3').find('dl', 'rank_list_box').find('h2').text.strip()
23 lt3 = [resault]
24 resault = soup.find('div', 'rank_list table rankstyle3').find_all('div', 'rank_txt')[:9]
25 for m in resault:
26     lt3.append(m.text.strip())
27
28 ranking = set(lt1 + lt2 + lt3)
29 ranking.remove('做工的人電影版')
30
31 print(ranking)

```

{'做工的人 電影版', '疫起', '龍與地下城：盜賊榮耀', '靈魂伴侶', '金牌拳手3', '驚聲尖叫6', '魔女宅急便', '捍衛任務4', '蟻人與黃蜂女：量子狂熱', '}

```

1 movie_infos = []
2 n = 0
3
4 for i in range(MOVIE_TOTAL_NUM):
5     link = LINK + str(i + 1)
6     links = []
7     labels = []
8
9     response = requests.get(link)
10    soup = BeautifulSoup(response.text, "lxml")

```

```

11
12     try:
13         resault = soup.find('div', 'movie_intro_info_r')
14         cname = resault.find('h1').text.strip()
15         ename = resault.find('h3').text.strip()
16         released_date = resault.find('span', class_=None).text[5:]
17
18         resault = soup.find('div', 'level_name_box').find_all('div', 'level_name')
19         for label in resault:
20             labels.append(label.text.strip())
21
22         resault = soup.find('div', 'gray_infobox_inner')
23         intro = str(resault.select_one('span').text).strip().replace('\n', '').replace('\r', '').replace(' ', ' ')
24
25         try:
26             resault = soup.find('ul', 'maylike_list_starlist').find_all('li', 'gabtn')
27             regex = re.compile(r'[0-9]+$')
28             for movie in resault:
29                 links.append(int(regex.findall(movie.find('a')['href'])[0]))
30         except:
31             pass
32     except:
33         continue
34
35     movie_info = MovieInfo(doc_id=i+1, cname=cname, ename=ename, label=labels, intro=intro, released_date=released_date)
36     movie_infos.append(movie_info)
37     n += 1
38
39 print(n)

```

11450

```

1 G = nx.DiGraph()
2 for movie_info in movie_infos:
3     for link_target in movie_info.links:
4         G.add_edge(movie_info.doc_id, link_target)
5 pagerank_list = nx.pagerank(G, alpha=1)
6 sorted_pagerank_list = sorted(pagerank_list.items(), reverse=True)
7 for movie_info in movie_infos:
8     item = sorted_pagerank_list.pop()
9     movie_info.pagerank = item[1]
10
11 movie = [movie_info.generate_dict() for movie_info in movie_infos]
12
13 with open("hw2.json", "w") as f:
14     json.dump(movie, f, indent=4)
15

```

```

1 inverted_index = {}
2 for movie in movie_infos:
3     seg = jieba.cut_for_search(movie.cname)
4     for word in seg:
5         if word in inverted_index.keys():
6             if movie.doc_id in inverted_index[word]:
7                 continue
8             inverted_index[word].append(movie.doc_id)
9         else:
10             inverted_index[word] = [movie.doc_id]
11
12     seg = jieba.cut_for_search(movie.intro)
13     for word in seg:
14         if word in inverted_index.keys():
15             if movie.doc_id in inverted_index[word]:
16                 continue
17             inverted_index[word].append(movie.doc_id)
18         else:
19             inverted_index[word] = [movie.doc_id]
20
21 with open("inverted.json", "w") as f:
22     json.dump(inverted_index, f, indent=4)
23

```

```
1 with open("inverted.json") as f:
2     inverted_index = json.load(f)
3
4 search_engine = SearchEngine(inverted_index, "hw2.json")
5 search_engine.query("捍衛任務")
6
```

您的搜尋結果:

共41筆, 符合"捍衛任務"

13573 (1.5808705135843418e-06):

Chinese Name: 殺戮基地

English Name: Black Site

Intro: ★《怒火邊界》《捍衛任務》製片打造黑暗系火爆動作鉅片★《不可能的任務系列》動作女星蜜雪兒摩納漢挑大樑主演★《魔鬼終結者：創世契機》男

14653 (1.0647066239985156e-06):

Chinese Name: 捍衛任務4

English Name: John Wick: Chapter 4

Intro: ★台灣搶先全球上映, IMAX、Dolby Cinema版本同步上映★系列全球賣座近6億《捍衛任務系列》原班人馬打造最新史詩篇章★「葉問」對決「殺神」

12776 (8.687221522247683e-07):

Chinese Name: 救命緝約

English Name: The Contractor

Intro: ★《捍衛任務》《天劫倒數》製作團隊最爽動作鉅獻★《神力女超人》克里斯潘恩《地獄》班佛斯特《恐懼大街》吉利安雅各布斯 大銀幕超激火拼★

14703 (8.292053684360614e-07):

Chinese Name: 驚爆點

English Name: Point Break

Intro: ★《危機倒數》奧斯卡金獎導演凱薩琳畢格羅執導★《捍衛任務系列》基努李維27歲嶄露頭角成名作★《鐵達尼號》《阿凡達：水之道》金獎名導詹姆

12899 (2.0322458521049444e-07):

Chinese Name: 非甜蜜生活

English Name: Three Floors

Intro: ★坎城影展競賽片 首映全場鼓掌十五分鐘★金棕櫚大師 南尼莫瑞提 睽違六年 深情力作★以色列動人小說改編 義大利金獎演員同台飆戲三個樓層、

14759 (1.345739939951112e-07):

Chinese Name: 殺手迴戰

English Name: Assassin Club

Intro: ★《金牌特務：金士曼起源》《天能》團隊傾力製作 火力全開★《不可能的任務7》《捍衛任務2》動作團隊玩命演出★《玩命快遞：肆意橫行》導演爽度

11346 (1.1019146165282226e-07):

Chinese Name: 一級任務

English Name: Voyagers

Intro: ★壓抑慾望被解封, 無重力多人廝殺混戰★《鋼鐵英雄》《捍衛任務》億萬製片打造, 超越想像嶄新烏托邦★《藥命效應》《分歧者》導演尼爾伯格 探

10821 (1.0303354377587279e-07):

Chinese Name: 追殺艾娃

English Name: Ava

Intro: ★《神鬼認證》《玩命關頭》幕後團隊打造暑期最強動作鉅片★《牠：第二章》潔西卡雀絲坦監製主演, 魅力打造女力版《捍衛任務》★《紳士追殺令》

11556 (9.944835477592685e-08):

Chinese Name: 豬殺令

English Name: Pig

Intro: ★《仲夏魘》《怪獸與葛林戴華德的罪行》製片團隊打造豬豬版《捍衛任務》! ★國際知名「爛番茄」影評網96%超高好評推薦★一切只為一頭豬! 尼

7820 (8.808257703044979e-08):

Chinese Name: 死侍2

English Name: Deadpool 2

Intro: 《死侍2》不只包括從「死侍」本人萊恩雷諾斯再次親手擔綱本片監製與男主角的重要角色, 更是首度好萊塢動作片名導《捍衛任務》大衛雷奇合作。而首

4825 (6.199851197712683e-08):

Chinese Name: 第七傳人

English Name: Seventh Son

Intro: 在傳說和巫術相互激盪的魔法時代, 全世界僅存的聖殿騎士葛瑞戈里(奧斯卡得主傑夫布里吉 飾, 《真實的勇氣》、《鋼鐵人》)周遊各地, 尋找擁有神

