

Hw#2, NLP@CGU Spring 2023 LINK: paste your link here

<https://colab.research.google.com/drive/1fN-LoT3F7N6Yv0Oy7v4nTfusl5E30ys6?usp=sharing>

Student ID: B0928001 Name: 賴霆瑞

```
import requests
from bs4 import BeautifulSoup
import networkx as nx

# specify the base URL and the range of pages you want to scrape
base_url = "https://movies.yahoo.com.tw/movieinfo_main/"
start_page = 1
end_page = 10

# define a Graph object
G = nx.Graph()

# loop through each page and scrape the movie information
for page in range(start_page, end_page+1):
    url = base_url + str(page)
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")
    if (soup.find("div", class_="movie_intro_info") == None):
        continue
    # add the URL to the Graph
    G.add_node(url)
    # find the movie information on the page
    doc_id = page
    cname = soup.find("div", class_="movie_intro_info").find("h1").text.strip()
    ename = soup.find("div", class_="movie_intro_info_r").find("h3").text.strip()
    label = soup.find("div", class_="level_name").text.strip()
    intro = soup.find("span", id="story").text.strip().replace('\n\n', '')
    released_date = soup.find("div", class_="movie_intro_info_r").find_all("span")[0].text.strip().replace('上映日期: ', '')

    # get all the links on the page and add them to the Graph
    links = []
    for link in soup.find_all("a"):
        href = link.get("href")
        text = link.text.strip()
        if href and text:
            links.append({"href": href, "text": text})
            G.add_edge(url, href)

    # print the movie information with the calculated PageRank
    pagerank = nx.pagerank(G)
    print({
        "doc_id": doc_id,
        "cname": cname,
        "ename": ename,
        "pagerank": pagerank[url], # use the URL as the key to get the PageRank
        "label": label,
        "intro": intro,
        "released_date": released_date,
        "links": links
    })

{'doc_id': 1, 'cname': '一世狂野', 'ename': 'Blow', 'pagerank': 0.46051496928893304, 'label': '劇情', 'intro': '喬治戎格一生都在追求所謂的美國夢', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 2, 'cname': '玩命關頭', 'ename': 'The Fast and the Furious', 'pagerank': 0.2410828215489779, 'label': '動作', 'intro': '唐米尼杜洛托', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 3, 'cname': '戰雲密佈', 'ename': 'Storm Catcher', 'pagerank': 0.14111163386770906, 'label': '動作', 'intro': '美國空軍最高機密的隱形', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 4, 'cname': '騎士風雲錄', 'ename': 'A Knight's Tale', 'pagerank': 0.11540232994260224, 'label': '動作', 'intro': '14世紀中古時期的社', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 5, 'cname': '金法尤物', 'ename': 'Legally Blonde', 'pagerank': 0.09584932468224716, 'label': '喜劇', 'intro': '在【歡樂谷】、【危險性', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 6, 'cname': '瘋狂世界', 'ename': 'Rat Race', 'pagerank': 0.07773888216344625, 'label': '冒險', 'intro': '薇拉貝克(琥碧戈珀飾)來到賭', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 7, 'cname': '震撼教育', 'ename': 'Training Day', 'pagerank': 0.07119912160242191, 'label': '劇情', 'intro': '在【震撼教育】中, 丹佐扮', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 8, 'cname': '神鬼第六感', 'ename': 'The Others', 'pagerank': 0.055190541423527265, 'label': '劇情', 'intro': '一個名叫葛瑞絲(妮可基嫻', 'released_date': '1997年剛搬家的星賢在信箱'}
{'doc_id': 10, 'cname': '觸不到的戀人', 'ename': 'il mare', 'pagerank': 0.04979482266247514, 'label': '奇幻', 'intro': '1997年剛搬家的星賢在信箱'}
```

```
### 爬取Yahoo! 電影資料
import requests
from bs4 import BeautifulSoup
import networkx as nx
import json

# specify the base URL and the range of pages you want to scrape
base_url = "https://movies.yahoo.com.tw/movieinfo_main/"
start_page = 1
end_page = 15067
```

```

# define a Graph object
G = nx.Graph()

# loop through each page and scrape the movie information
for page in range(start_page, end_page+1):
    url = base_url + str(page)
    response = requests.get(url)
    soup = BeautifulSoup(response.content, "html.parser")
    if (soup.find("div", class_="movie_intro_info") == None):
        continue
    # add the URL to the Graph
    G.add_node(url)
    # find the movie information on the page
    doc_id = page
    cname = soup.find("div", class_="movie_intro_info").find("h1").text.strip()
    ename = soup.find("div", class_="movie_intro_info_r").find("h3").text.strip()
    label = soup.find("div", class_="level_name").text.strip()
    intro = soup.find("span", id="story").text.strip().replace('\n\n', '')
    released_date = soup.find("div", class_="movie_intro_info_r").find_all("span")[0].text.strip().replace('上映日期: ', '')

    # get all the links on the page and add them to the Graph
    links = []
    for link in soup.find_all("a"):
        href = link.get("href")
        text = link.text.strip()
        if href and text:
            links.append({"href": href, "text": text})
            G.add_edge(url, href)

    # calculate the PageRank for the page
    pagerank = nx.pagerank(G)

    # create a dictionary for the movie information
    movie_info = {
        "doc_id": doc_id,
        "cname": cname,
        "ename": ename,
        "pagerank": pagerank[url], # use the URL as the key to get the PageRank
        "label": label,
        "intro": intro,
        "released_date": released_date,
        "links": [link["href"] for link in links] # store only the hrefs of the links
    }

    # append the movie information to the existing JSON file
    with open("movies.json", "a", encoding="utf-8") as f:
        json.dump(movie_info, f, ensure_ascii=False)
        f.write("\n") # add a new line at the end of each JSON object to separate them

### 中文分詞後，建立 Inverted Index
import json
import jieba
from collections import defaultdict

# load the movie data from the JSON file
with open('movies.json', 'r', encoding='utf-8') as f:
    movies = [json.loads(line) for line in f]

# define a tokenizer function using jieba
def tokenize(text):
    return list(jieba.cut(text))

# build the inverted index
inverted_index = defaultdict(list)
for movie in movies:
    # tokenize the movie intro, cname, ename, and label fields
    tokens = tokenize(movie['intro'] + movie['cname'] + movie['ename'] + movie['label'])
    # add each token to the inverted index along with the movie ID
    for token in tokens:
        inverted_index[token].append(movie['doc_id'])

# save the inverted index to a JSON file
with open('inverted_index.json', 'w', encoding='utf-8') as f:
    json.dump(inverted_index, f, ensure_ascii=False)
    f.write("\n") # add a new line at the end of each JSON object to separate them

Building prefix dict from the default dictionary ...
DEBUG:jieba:Building prefix dict from the default dictionary ...
Dumping model to file cache /tmp/jieba.cache
DEBUG:jieba:Dumping model to file cache /tmp/jieba.cache
Loading model cost 1.005 seconds.

```

```
DEBUG:jieba>Loading model cost 1.005 seconds.
Prefix dict has been built successfully.
DEBUG:jieba:Prefix dict has been built successfully.
```

```
### 利用 PageRank 演算法來排序
import json

# read in the movie information from the JSON file
movies = []
with open("movies.json", "r", encoding="utf-8") as f:
    for line in f:
        movie_info = json.loads(line)
        movies.append(movie_info)

# sort the movies based on their PageRank scores
movies_sorted = sorted(movies, key=lambda x: x["pagerank"], reverse=True)

# print out the top 10 movies
for i, movie in enumerate(movies_sorted[:50]):
    print(f"{i+1}. {movie['cname']} ({movie['released_date']}) - PageRank score: {movie['pagerank']:.4f}")

1. 黑狗來了 (2004-03-27) - PageRank score: 0.4606
2. 我和吸血鬼有份合約 (2001-11-24) - PageRank score: 0.4606
3. 北京樂與路 (2001-11-10) - PageRank score: 0.4605
4. 幸福騙局 (2018-11-30) - PageRank score: 0.4605
5. 麻雀變公主 (2001-09-01) - PageRank score: 0.4605
6. 一世狂野 (2001-10-12) - PageRank score: 0.4605
7. 鬼訊號 (2005-01-14) - PageRank score: 0.4605
8. 鬼仔 (2013-11-29) - PageRank score: 0.4605
9. 英倫對決 (2017-10-13) - PageRank score: 0.4605
10. 野鴨變鳳凰：冰上逆襲 (2021) ((2021)) - PageRank score: 0.4605
11. 泰德拉索：錯棚教練趣事多 (2021) ((2021)) - PageRank score: 0.4603
12. 野鴨變鳳凰：冰上逆襲 (2022) ((2022)) - PageRank score: 0.2485
13. 實尾島風雲 (2004-04-02) - PageRank score: 0.2422
14. 玩命關頭 (2001-10-13) - PageRank score: 0.2411
15. 關鍵琴聲 (2014-01-17) - PageRank score: 0.2394
16. 千禧曼波 (2001-11-17) - PageRank score: 0.2363
17. BJ單身日記 (2001) (2001-10-27) - PageRank score: 0.2343
18. 泰德拉索：錯棚教練趣事多 (2020) ((2020)) - PageRank score: 0.2319
19. 血肉森林 (2004-12-24) - PageRank score: 0.2288
20. 這旅程使命必達 (2018-12-07) - PageRank score: 0.2288
21. 徐自強的練習題 (2017-08-11) - PageRank score: 0.2216
22. 晚孃 (2001-11-03) - PageRank score: 0.2209
23. 屏息 (2013-12-13) - PageRank score: 0.1764
24. 變身特務 (2020-01-22) - PageRank score: 0.1717
25. 見鬼2 (2004-03-26) - PageRank score: 0.1602
26. 青禾男高 (2017-07-21) - PageRank score: 0.1597
27. 救世主 (2001-12-07) - PageRank score: 0.1570
28. 人間有情天 (2001-11-17) - PageRank score: 0.1530
29. 貝蒂費雪的世界 (2004-12-10) - PageRank score: 0.1516
30. 追夢高手 (2001-11-30) - PageRank score: 0.1447
31. 心碎效應 (2022-09-23) - PageRank score: 0.1441
32. 戰雲密佈 (2001-10-13) - PageRank score: 0.1411
33. 御史與祚怡 (2021) ((2021)) - PageRank score: 0.1395
34. 功夫 (2004-12-24) - PageRank score: 0.1244
35. 老師君主 (2018-11-16) - PageRank score: 0.1228
36. 美國派2 (2001-11-10) - PageRank score: 0.1211
37. 行動代號：狼狩獵 (2022-09-30) - PageRank score: 0.1204
38. 騎士風雲錄 (2001-10-19) - PageRank score: 0.1154
39. 人性污點 (2004-04-02) - PageRank score: 0.1136
40. 如蝶翩翩 (2021) ((2021)) - PageRank score: 0.1118
41. 鬼計神偷 (2001-10-06) - PageRank score: 0.1098
42. 愛貓之城 (2017-08-04) - PageRank score: 0.1080
43. 屍房宴 (2014-01-17) - PageRank score: 0.1072
44. 大發不動產 (2021) ((2021)) - PageRank score: 0.1043
45. 哥兒們 (2001) (2001-11-17) - PageRank score: 0.1042
46. 靈魂的重量 (2004-04-16) - PageRank score: 0.0989
47. 艾蜜莉的異想世界 (2002-03-08) - PageRank score: 0.0989
48. 金法尤物 (2001-10-19) - PageRank score: 0.0958
49. 少林足球 (2001-08-24) - PageRank score: 0.0944
50. 童話·世界 (2022-10-07) - PageRank score: 0.0930

### 定義搜尋Function
import json

def search_movies(query):
    # load the movie info from the JSON file
    with open("movies.json", "r", encoding="utf-8") as f:
```



```
# load the movie info from the JSON file
with open("movies.json", "r", encoding="utf-8") as f:
    movies = [json.loads(line) for line in f]

# create a dictionary of {doc_id: PageRank value} pairs
pageranks = {movie["doc_id"]: movie["pagerank"] for movie in movies}

# create a list of movies that match the search query
matches = [movie for movie in movies if any([query.lower() in str(value).lower() for key, value in movie.items()])]

# sort the matches by their PageRank values
sorted_matches = sorted(matches, key=lambda x: pageranks.get(x["doc_id"], 0), reverse=True)

# calculate precision and recall
relevant_docs = set([movie['doc_id'] for movie in movies if query.lower() in str(movie).lower()])
retrieved_docs = set([movie['doc_id'] for movie in matches])
true_positives = relevant_docs.intersection(retrieved_docs)
precision = len(true_positives) / len(matches) if len(matches) > 0 else 0
recall = len(true_positives) / len(relevant_docs) if len(relevant_docs) > 0 else 0

# print the search results and evaluation metrics
print("您的搜尋結果 (Sorting by PageRank Value): ")
print(f"共 {len(sorted_matches)} 筆, 符合 '{query}' - - - 共 indexing {len(movies)} 筆電影資料")
for movie in sorted_matches:
    intro = movie.get("intro", "")
    print(f"{movie['doc_id']} ({pageranks.get(movie['doc_id'], 0)}): {movie['cname']} ({movie['ename']}) - {intro}")

print(f"\nPrecision: {precision:.0%}")
print(f"Recall: {recall:.0%}")
```

```
### 輸入搜尋關鍵字, 任何字都可以
search_movies("Blow")
```

▲ 瑞秋·懷茲、吳姓下只路琪，兩人又生聯才無感。又言奈凱薩琳，飯上徽拍抓封級四化星，加入聯官國經又部隊，協助里廷的監官上TF。她的任務是訓練級四化星4766 (0.00017984344657567337)：戀戀海灣 (As the Winds Blow) - 在海邊出生，幼時隨著母親北上定居，米雪早已忘記自己體內的海洋基因。在母親的命令下

Precision: 100%