

▼ Lab#3, NLP@CGU Spring 2023

This is due on 2023/03/20 16:00, commit to your github as a PDF (lab3.pdf) (File>Print>Save as PDF).

IMPORTANT: After copying this notebook to your Google Drive, please paste a link to it below. To get a publicly-accessible link, hit the *Share* button at the top right, then click "Get shareable link" and copy over the result. If you fail to do this, you will receive no credit for this lab!

LINK: paste your link here

<https://colab.research.google.com/drive/12kAOrHflnFHbTjrvlOMGK7pmTdFNwKJr>

Student ID: B0928005

Name: 湯嘉為

▼ Question 1 (100 points)

Implementing Yahoo Movies Crawler.

1. Design a Yahoo! Movie Crawler.
2. Crawl all the movie information listed in movie_intheaters page
3. The more movie data crawled, the higher the score

按兩下 (或按 Enter 鍵) 即可編輯

```
import requests
import re
from bs4 import BeautifulSoup

Y_MOVIE_URL = "https://movies.yahoo.com.tw/movie_intheaters.html"

# YOUR CODE HERE!
# IMPLEMENTING YAHOO MOVIES CRAWLER

class MovieCrawler(object):

    def __init__(self):
        self.headers = {
            'User-Agent': 'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36 (KHTML, like Gecko) Chrome/58.0.3029.110'
        }
        self.movies = []

    def get_movies(self, page_url):
        res = requests.get(page_url, headers=self.headers)
        soup = BeautifulSoup(res.text, 'html.parser')
        movie_list = soup.find_all('div', class_='release_info')
        for movie in movie_list:
            # Get the Chinese name of the movie.
            ch_name = movie.find('div', class_='release_movie_name').a.text.strip()
            # Get the English name of the movie.
            en_name = movie.find('div', class_='release_movie_name').find('div', class_='en').a.text.strip()
            # Get the URL of the movie detail page.
            movie_url = movie.find('div', class_='release_movie_name').a['href']
            # Get the URL of the movie detail page.
            release_date = movie.find('div', class_='release_movie_time').text.strip()
            release_date = release_date.split(' ')[-1]
            ## Get the introduction of the movie.
            intro = movie.find('div', class_='release_text').text.strip().replace('\n', '').replace('r', '')

            # Store the movie information in a dictionary and append it to the movies list.
            movie_dict = {
                'ch_name': ch_name,
                'en_name': en_name,
                'movie_url': movie_url,
                'release_date': release_date,
                'intro': intro
            }
            self.movies.append(movie_dict)

        # Find the link to the next page
        next_page = soup.find('a', rel='next')
        if next_page:
            # Construct the link to the next page and recursively call the get_movies() function.
            next_page_url = next_page['href']
            self.get_movies(next_page_url)
```

```

        return self.movies

# # DO NOT MODIFY THE VARIABLES
crawler = MovieCrawler()
movies = crawler.get_movies(Y_MOVIE_URL)

# # THE RESULTS : AS THE FOLLOWING SECTION
# # {'ch_name', 'en_name', 'movie_url', 'release_date', 'intro'}
print(len(movies))
print(*movies, sep="\n")

```

者，以及500多首令人難忘的電影配樂的作曲家，本片則堪稱是他最完整肖像電影。奧斯卡導演托納多雷片中訪問了包括與莫利克奈合作過的王家衛、昆汀塔倫蒂諾（《利器》）以及艾美獎得主瑪果麥丁達爾（諜報影集《冷戰諜夢》）和艾美獎得主雷里歐塔（《四海好傢伙》）。《熊蓋毒》一片的導演是伊莉莎白班克斯（「是以此為理由離開對方呢？」）

阿曼（《Second Chances》）飾演佩卓潘納；瑪莎米蘭絲（影集《白線》）飾演蘿莎范斯奎茲；庫柏安德魯斯（影集《陰屍路》）飾演維克多范斯奎茲；以及吉蒙

惡行。第一個秘密過去，哇你真夠扯的令人髮指……第二個秘密過去，原來你不只壞而且連靈魂都壞光光了吧……第三個秘密……等等！這第三個秘密不及格吧？

八竿子都打不著關係的這兩個家庭，也迎來天翻地覆的變化……《玩具當家》翻拍自1976年的法國社會喜劇電影《The Toy》，由法國喜劇泰斗——賈梅德布茲飾演

不同往昔。’}

。【關於電影】《愛麗絲與藏六》今井哲也同名科幻漫畫改編！轟動國際各大影展的話題動畫電影《我們的黎明》，改編自漫畫家今井哲也的同名科幻漫畫作品，族綁架…直到他遇見女科學家伊娃（莉蒂西亞杜希飾），她完全改變了文森的价值觀，現在的他，將作出人生中最勇敢的決定…」

最佳男主角的許光漢、金馬最佳男配角的林柏宏、台北電影節影后王淨，三位人氣及實力兼具的演員主演，集結蔡振南、庾宗華、馬念先等實力派演員共同演出，山林的記錄者，看見前人的戮力以赴。」}

改寫紀錄的女性桂冠指揮家。身為女性，瑪琳艾索普所取得的非凡成就來並不容易。幼時隨著父親聆賞李奧納多•伯恩斯坦 (Leonard Bernstein) 演出的震撼體
 名★ 2023 棕櫚泉國際影展 聚光燈獎 - 布蘭登費雪一名英語老師查理 (布蘭登費雪 飾) 中年出櫃後拋棄家庭與愛人私奔，如今他因為愛人過世而罹患暴戾症，同

