

▼ 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/1_Gw6MaB7HHa9Tor46aLUF9jPmyeoNe3n?usp=sharing

Student ID: B0928001 **Name:** 賴霆瑞

▼ 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

```
import requests
from bs4 import BeautifulSoup

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

class MovieCrawler(object):
    def __init__(self):
        pass

    def get_movies(self, page_url):
        # 創建一個空列表，用於存儲電影資訊
        movies = []

        # 使用requests模組發送GET請求獲取網頁內容
        resp = requests.get(page_url)

        # 解析HTML內容
        soup = BeautifulSoup(resp.content, 'html.parser')

        # 查找所有電影資訊的HTML元素
        movie_items = soup.find_all('div', {'class': 'release_info'})

        # 循環處理每個電影資訊元素，提取所需的資訊
        for item in movie_items:
            movie = {}
            # 獲取電影中文名
            movie['ch_name'] = item.find('div', {'class': 'release_movie_name'}).a.text.strip()
            # 獲取電影英文名
            movie['en_name'] = item.find('div', {'class': 'en'}).a.text.strip()
            # 獲取電影詳細頁面的URL
            movie['movie_url'] = item.find('div', {'class': 'release_movie_name'}).a['href']
            # 獲取電影上映日期
            movie['release_date'] = item.find('div', {'class': 'release_movie_time'}).text.strip().replace('上映日期: \n', '')
            # 獲取電影簡介
            movie['intro'] = item.find('div', {'class': 'release_text'}).text.strip().replace('\r\n', '').replace('\n', '')

            # 將電影資訊添加到列表中
            movies.append(movie)

        return movies

# 創建MovieCrawler對象，獲取前3頁電影資訊
crawler = MovieCrawler()

for page in range(1, 9):
    page_url = f"{Y_MOVIE_URL}?page={page}"
    movies = crawler.get_movies(page_url)
    # 輸出當前頁數和電影數量
    print(f"Page {page} - Total movies: {len(movies)}")
    # 輸出電影資訊
    print(*movies, sep="\n")

# 創建MovieCrawler對象，獲取電影資訊
# crawler = MovieCrawler()
```


✓ 4 秒 完成時間: 下午3:58

● ×