

▼ 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/1btirJ2gncleKojDu5-BYQYdl7ir8TTya?usp=sharing>

Student ID:B092800 **Name:**郭玉俊

▼ 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
import csv

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

    def get_movies(self, page_url):
        movies = []
        page = 1
        while True:
            response = requests.get(url=page_url, params={'page': page})
            soup = BeautifulSoup(response.text, 'lxml')
            info_items = soup.find_all('div', 'release_info')
            if not info_items: # 如果找不到任何電影, 就退出循環。
                break
            for item in info_items:
                movie = {}
                movie['ch_name'] = item.find('div', 'release_movie_name').a.text.strip()
                movie['en_name'] = item.find('div', 'en').a.text.strip()
                movie['movie_url'] = item.find('div', 'release_movie_name').a['href'].strip()
                movie['release_date'] = item.find('div', 'release_movie_time').text.split(':')[1].strip()
                movie['intro'] = item.find('div', 'release_text').span.text.strip()
                movies.append(movie)
            page += 1
        return 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")
```

