

▼ 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/1s4JEY5DijL2doRCTrh-40mbr_9kt3M2n?usp=sharing

Student ID:B0928013

Name:吳佳恩

▼ 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"
response = requests.get(url=Y_MOVIE_URL)

soup = BeautifulSoup(response.text, 'lxml')
# YOUR CODE HERE!
# IMPLEMENTIG YAHOO MOVIES CRAWLER

class MovieCrawler(object):

    def __init__(self):
        info_items = soup.find_all('div', 'release_info')
    def get_movies(self, page_url):
        for item in info_items:
            chinese_name = item.find('div', 'release_movie_name').a.text.strip()
            english_name = item.find('div', 'en').a.text.strip()
            release_time = item.find('div', 'release_movie_time').text.split(':')[1].strip()
            moive_url = item.find('a', class_="gabtn")['href']
            movie_text = item.find('span', 'jq_text_overflow_180 jq_text_overflow_href_list').text

            #print('chname: {} enname: ({} ) 上映日: {} moive_url: {} movie_text {}'.format(chinese_name, english_name, release_time, moive_ur

get_movies(self, page_url)
# # 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")

-----
TypeError                                 Traceback (most recent call last)
<ipython-input-17-43bd994892c7> in <module>
    33 # # THE RESULTS : AS THE FOLLOWING SECTION
    34 # # {'ch_name', 'en_name', 'movie_url', 'release_date', 'intro'}
--> 35 print(len(movies))
    36 print(*movies, sep="\n")

TypeError: object of type 'NoneType' has no len()
```

SEARCH STACK OVERFLOW

[illegible]

✓ 0 秒 完成時間: 下午3:52

● ×