# Assignment 1: Web Scraping & Data Analysis

Oct 31, 2021

In this assignment, you should work with data from <https://maoyan.com/board/4> (Top 100 Movie)

Maoyan Movies is a favorite viewing platform for domestic audiences and provides you with online ticket purchase services. At the same time, Maoyan Movies also provides you with movie trailers, box office queries, movie rankings, film and television information and other information.

Everyone is interested in bad movies, how can there be so many bad movies? So, it is very important to scrape high-quality movies on various websites. In this project, we will use the requests and regular expressions we have learned before **to scrape the TOP100 movies from Maoyan Movies.**

If you CANNOT read Chinese, please crawl other websites (i.e., https://brickset.com/sets/year-2016). Then report should be written based on your collected data. <u>On the report, please specify which URL you are working with.</u>

*<u>Task1</u>.* I want you to **scrape the latest 100 movies** from the website and save result into 'your_name+id.csv'. This file should contain the data with the following columns: (40pt)

- Title
- Name of director
- Name of actors
- Rating
- Cumulative income
- Duration
- Type

You are free to explore data **with more properties** if needed.

*<u>Task2</u>.* I want you to do a data analysis on the data. What do you think is interesting about this data? Tell a story **within 2 pages** (excluding figures) about some interesting thing you have discovered by looking at the data. (60pt)

For example, you might consider whether the director has an impact on movie box office revenue, or whether some directors only focus on making certain type of moves. Another thing you might consider is whether there is a relationship between the box office revenue and the user's ratings on the movie.

The assignment code that runs on its own (Web crawling + Data analysis) should be handed in using a Jupter Notebook file.

**Submission Checklist:**

| Yes/No | Items |
|---|---|
| | Assignment code |
| | your_name+id.csv |
| | 2 pages Repoort |