
839 Project Stage 2: Extracting Movie data from IMDB and MovieNumbers

Daniel K. Griffin
Department of Computer Science
University of Wisconsin Madison
dgriffin5@wisc.edu

Yudhister Satija
Department of Computer Science
University of Wisconsin Madison
ysatija@wisc.edu

Mitali Rawat
Department of Computer Science
University of Wisconsin Madison
mitali.rawat@cs.wisc.edu

1 Introduction

For the Data Science 839 project stage 2, we were asked to develop rule based extractors to extract structured data from two overlapping sources of our own choosing, and to normalize both data sets to have the same schema. We decided to extract movies from the IMDB movies website, and the Movie Numbers website using a combination of standard python scripts, and the Scrapy framework for building web crawlers. A quick description of how we went about the task is given at a high level below, and the commands for how to run our web crawlers is provided in our root github README.md web page.

2 The Quick List

This section provides a quick reference sheet of our accomplishments, matching the list of what the .pdf file should contain on the course project web page. This section is meant to provide a quick synopsis of the required information to make evaluation easier.

1. **The Web data sources:** We used the following 2 sources:

- IMDB movies list from <http://www.imdb.com/list/ls032600534>
- Movie Numbers list from <https://www.the-numbers.com/movies/#tab=letter>

2. **How we extracted structured data:**

- **Imdb extraction:** Extraction from the Imdb movie database website was a two stage process. The first stage was to parse data from the website using a Scrapy based web crawler, parse the data from the page using css based DOM selectors, and generate a .json data store of the parsed information. This first stage was fairly straightforward since the data was relatively flat, and there were no nested links that needed to be parsed. The second stage was to convert the .json data store into a .csv database type file.
- **Movie numbers extraction:** Extraction from the Movie Numbers website followed the same two stage process as the Imdb extraction. The main difference was that extraction from this data source was more complicated, as multiple nested links needed to be traversed at multiple pages to extract all of the necessary movies and movie information. Also, our extractor for information from these pages relied on xpath based selectors, rather than css based selectors.

- 30 3. **Type of entity:** Our entities are movies. We have movie names and information like release
31 date, directors and few cast member names. We also have MPAA rating of the movie and
32 some figures on the sales.
- 33 4. **The two tables:** The two tables are included in the csv files:
- 34 • "imdb.csv" for IMDB movie list
 - 35 • "thenumbers.csv" for Move Numbers movie list.
- 36 The tables have the following identical schema:
37 **(title, year, mpaa, runtime, genres, star_rating, metacore_rating, description, direc-**
38 **tor, stars, gross)**
- 39 The "stars" field refers to a few cast member names, separated by comma.
40 The "genres" field refers to the genres the movie is categorized under, separated by comma.
41 The "gross" field refers to sales numbers earned by the movie worldwide.
- 42 5. **Number of tuples:** We have obtained 4,291 tuples from the IMDB list. We have obtained
43 31,006 tuples from the Move Numbers' list.
- 44 6. **Open-source tools:**
- 45 • **Python Scrapy:** Scrapy is an open-source python framework for quickly developing
46 complex website scraping web spiders. Scrapy was designed to be easily integrated into
47 a structured data parsing pipeline, and contains many integrated parsing modules for
48 logging information, generating data statistics, generating emails given parsing events,
49 and complex modules for parsing structured data using xpath and css DOM based
50 selectors, web page link extraction tools, and support for integrating more complicated
51 post-processing mechanisms for parsed data. We used Scrapy to develop web crawling
52 spiders for both the imdb and movie numbers web sites, and use both the integrated
53 xpath based and css DOM based selectors for pulling information from the structured
54 web pages.

55 3 The Data Sources

56 We extracted movie information from the IMDB and Movie Numbers websites. IMDB is a general
57 movie website that stores information about movies including categories of movies, related movies,
58 movie ratings, movie descriptions, movie critiques, and many other pieces of information. For
59 our purposes, we found a fairly substantial flat list of movies (about 4000 movies) structured in a
60 browseable format for 100 movies per page, with links at the bottom of each page to move from
61 one page to the next. Movie Numbers is another website that contains movie information, and holds
62 roughly the same kinds of information as IMDB. The Movie Numbers website has an index of movies
63 that allows a user to search for movies alphabetically based on movie name, and allows access to
64 movies based off of a nested weblink indexing structure.

65 4 The Extractors

66 Both extractors are constructed of two smaller scripts. One script contains the code for a Scrapy
67 based web crawler, and the other script contains the code for calling the crawler, and parsing its
68 output into a .csv based file format. For the IMDB website, "IMDBSpider.py" contains the web
69 crawler code, and "scrapeIMDB.py" contains the code that runs the IMDB spider and parses its
70 output into a .csv file format. For the Movie Numbers website, "NumbersSpider.py" contains the web
71 crawler code, and "scrapeMovieNumbers.py" contains the code that runs the Movie Numbers spider
72 and parses its output into a .csv file format.

73
74 The "IMDBSpider.py" script is a fairly simple script that retrieves a page of moves, and
75 uses css based DOM selectors to select each movie from the list of movies in the web page, and
76 parses information from each movie. The overall structure of the crawler is fairly simple, as no
77 nested links need to be traversed to get all of the movies from each web page. Also, since the web url
78 follows a simple naming convention for all pages, the web crawler explores all web pages in parallel
79 by expanding the number of every page for the web url, which speeds up parsing significantly. The
80 "scrapeIMDB.py" web script calls the "IMDBSpider.py" script using a python subprocess to generate
81 a .json file. When finished, the script reads the .json file, and uses python's built in json module to

82 manipulate the json data, and write the json data into a csv based format.

83

84 The Movie Numbers extractor works almost exactly the same as the IMDB based extractor.
85 The main difference however is in the complexity of the "NumbersSpider.py" web crawler. The
86 Movie Numbers website uses a nested web link indexing structure. So, the website web crawler
87 has to follow each nested link to find a movie information webpage before any information can be
88 extracted. This means the overall web crawler was more complicated than the IMDB based crawler,
89 as we were able to assume a flat indexing structure when parsing data from the IMDB based web
90 pages.