

Prediction of Movies Box Office Using Crawler

Mandar Choudhary, Kshiteej Naik, Shreyash Dange
Smt. Indira Gandhi College of Engineering, Navi Mumbai

Abstract - In this study, we apply data Web Crawler to generate interesting patterns for predicting box office performance of movies using data collected from multiple web sources including Bookmyshow, Nowrunning and the various movie database sources. The prediction is based on decision factors derived from a historical movie database. We label the prediction in three classes, Hit, Neutral and Flop.

Keywords: Web crawler, Wrapper, Link Finder, MySQL, Motion Pictures

I. INTRODUCTION

In this Digital age sites such as including Bookmyshow, Nowrunning and ratingdada have been used for sharing contents and reviews on all types of subjects of movies by millions of people on a daily basis. It is clear that businesses have a strong interest in tapping into these huge data sources to extract information that might improve their decision making process. For example, predictive models derived from so for successful movies may facilitate filmmakers making more profitable decisions. The topic of movies is of considerable interest in the social media user community. Research has been done to generate models for predicting the success of movies. The results are derived from multiple data source. These reviews are used to predict box-office sentiments of the movies. The Internet Movie Database (IMDb) is another popular data source for these types of studies. The prediction of movies performance were based on the film critics reviews where as in predictions were done based on regression or stochastic models on IMDb data. The approach combines social network analysis and automatic sentiment analysis.

Trends and real world events in movie business were predicted by weighing the forum posts on different rating sites.

II. PROBLEM STATEMENT

Existing systems do not target multiple sites for movie ratings instead predictions are done only on critics and reviews. Our system target multiple sites and gather data from various sources. Another problem is with the accuracy, our system is more accurate and success rate is high enough unlike other existing systems till date. One more factor with other existing systems is minimum range of operating depth and extensiveness. Our system uses complex operations and is more extensive in terms of operating depth by using complex crawling mechanisms.

And one important problem is use of K-means algorithm by existing systems which need input initiation value that can vary resulting in constant variation of output. This is avoided in our system by using DFS algorithm which maximizes the automation.

Proposed System

The proposed system is mostly automated, highly user friendly. Here we are going to implement a 'average movie rating' using crawling technique. The user has to enter the movie name. The average value will be fetched from different websites i.e., After searching the movie name the crawler will crawl all the links from the website respectively for each website. Then all the ratings from the websites will be fetched and by using averaging algorithm it will calculate the single average rating and display it to the user.

Advantages:

1. User satisfaction from search directed access to resources and easier browse ability (via maintenance and advancements of the Web resulting from analyses).
2. Reduced network traffic in document space resulting from search-directed access.
3. Effecting archiving/mirroring, and populating caches (to produce associated benefits).
4. Monitoring and informing users of changes to relevant areas of Web-space.
5. "Schooling" network traffic into localised neighbourhoods through having effected mirroring, archiving or caching.

III. PROPOSED SYSTEM WORKFLOW

A Web crawler, sometimes called a spider, is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing (*web spidering*). Web search engines and some other sites use Web crawling or spidering software to update their web content or indices of others sites' web content.

Web crawlers copy pages for processing by a search engine which indexes the downloaded pages so users can search more efficiently. Crawlers consume resources on visited systems and often visit sites without approval. Issues of schedule, load, and "politeness" come into play when large collections of pages are accessed. Mechanisms exist for public sites not wishing to be crawled to make this known to the crawling agent. For instance, including a robots.txt file can request bots to index only parts of a website, or nothing at all.

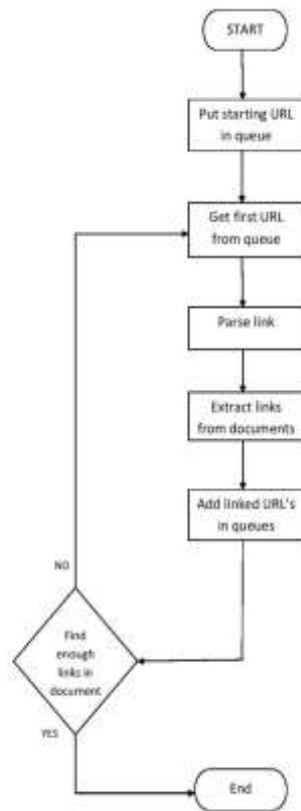


Fig. Basic workflow

IV. CONCEPTS USED Web Crawler

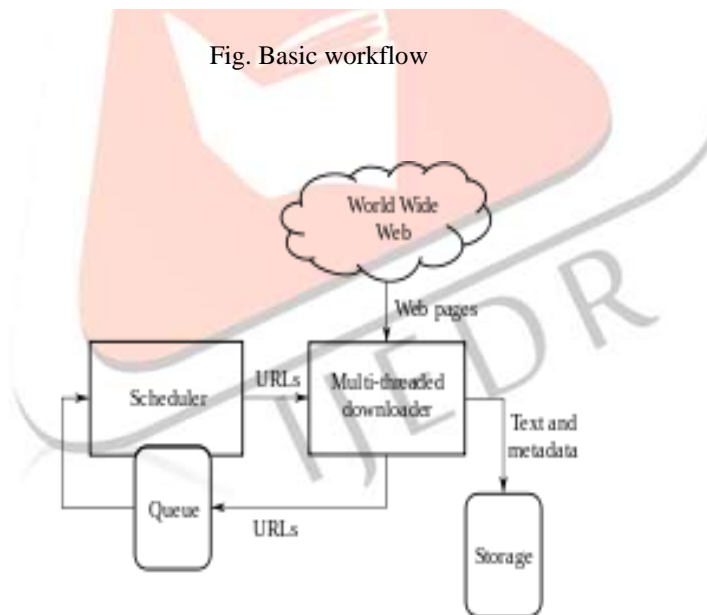


Fig. Architecture of a Web crawler

A Web crawler, sometimes called a spider, is an Internet bot that systematically browses the World Wide Web, typically for the purpose of Web indexing (*web spidering*).

Web search engines and some other sites use Web crawling or spidering software to update their web content or indices of others sites' web content. Web crawlers copy pages for processing by a search engine which indexes the downloaded pages so users can search more efficiently.

Crawlers consume resources on visited systems and often visit sites without approval. Issues of schedule, load, and "politeness" come into play when large collections of pages are accessed. Mechanisms exist for public sites not wishing to be crawled to make this known to the crawling agent. For instance, including a robots.txt file can request bots to index only parts of a website, or nothing at all.

The number of Internet pages is extremely large; even the largest crawlers fall short of making a complete index. For this reason, search engines struggled to give relevant search results in the early years of the World Wide Web, before 2000. Today relevant results are given almost instantly.

Crawlers can validate hyperlinks and HTML code. They can also be used for web scraping

Wrapper

When researchers translate their thoughts into codes, they usually have the needs to modify their existing functions to accommodate their new ideas, such as adding new arguments or a little more computations to the functions.

Wrapper functions can be used as an interface to adapt to the existing codes, so as to save you from modifying your codes back and forth.

The wrapper function typically performs some prologue and epilogue tasks like allocating and disposing resources, checking pre- and post-conditions, caching / recycling a result of a slow computation but otherwise it should be fully compatible with the wrapped function, so it can be used instead of it.

Beautiful Soup

Beautiful soup is a python package for parsing HTML and XML documents(including having malformed markup, i.e. non-closed tags, so named after tag soup). It creates a parse tree for parsed pages that can be used to extract data from HTML, which is useful for web scraping.

Link Finder

The hypertext reference, or href, attribute is used to specify a target or destination for the anchor element. It is most commonly used to define a URL where the anchor element should link to.

An href can do a lot more than just link to another website. It can be used to link directly to any element on a web page that has been assigned an id. It can be used to link to a resource using a protocol other than http. It can be used to run a script.

Database connectivity

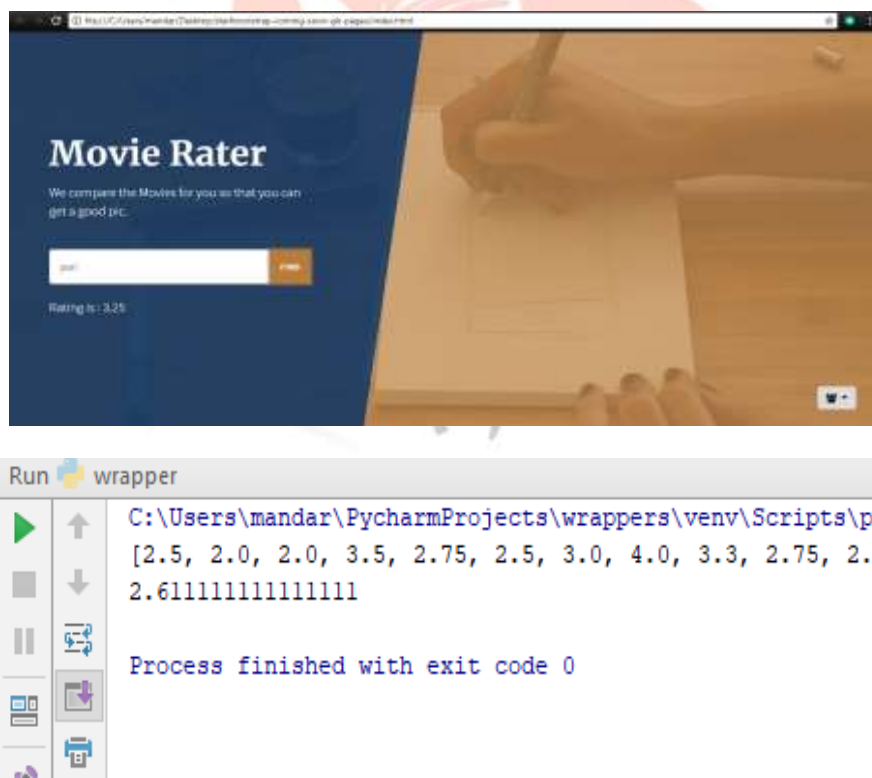
MySQL is well known as world's most widely used open-source database (back-end). It is most supportive database for PHP as PHP-MySQL is most frequently used open-source scripting database pair. The user-interface which WAMP, LAMP and XAMPP servers provide for MySQL is easiest and reduces our work to a large extent.

```
String db="ratings-mov";
```

```
Connectioncon=null;
```

```
String url = "jdbc:mysql://localhost/"+db; Class.forName("com.mysql.jdbc.Driver").newI
```

V. RESULT



VI. CONCLUSION

The model used to predict movie success is quite simple, but still powerful enough to make good predictions. Compared with other proposed methods from the challenge it is by far less sophisticated, though its strengths lie in its simplicity. The approach to viewer rating prediction is very straightforward and can be easily applied by anyone. It saves time for finding different sites, view the ratings, and then again access other sites for comparison. Thus it is an integrated approach of providing average movie ratings of various movies on a single platform on the move. This makes prediction of movie ratings very simple and effective.

Future Scope :

1. More Wrappers to target more websites.
2. Sentiment analysis to gather reviews based on comments.
3. Site to calculate multiple movies rating in a single iteration.

VII. REFERENCES

- [1] Sitaram Asur and Bernardo A. Huberman, "Predicting the Future with Social Media," <http://arxiv.org/abs/1003.5699>, March 2010.
- [2] Wenbin Zhang , Steven Skiena, Improving Movie Gross Prediction through News Analysis, Proceedings of the 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology, p.301-304, September 15-18, 2009 [doi>10.1109/WI-IAT.2009.53]
- [3]"YouTube", youtube.com, Inc., {Online}. Available: <http://www.youtube.com/yt/press/statistics.html>
- [4]"Coming soon movies", Inc., {Online}. Available: <http://www.comingsoon.net/movies.php>
- [5]"Movie box", Inc., {Online} Available: <http://themoviebox.net/view/trailers>
- [6]"Google Developers", Inc., {Online} Available: https://developers.google.com/youtube/2.0/developers_guide_protocol
- [7]"Internet movie database," IMDb.com, Inc., {Online}. Available: <http://imdb.com>. {Accessed 26 April 2013}.
- [8]"Twitter," Twitter.com, Inc., {Online}. Available: <http://twitter.com>. {Accessed 16 April 2013}.
- [9]"Sentimental Analysis", Inc. {Online}. Available: <http://www.cs.uic.edu/~liub/FBS/sentiment-analysis> {Accessed 23 March 2013}.
- [10] Mahesh Joshi , Dipanjan Das , Kevin Gimpel , Noah A. Smith, Movie reviews and revenues: an experiment in text regression, Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, p.293-296, June 02-04, 2010, Los Angeles, California
- [11] A. Chen, "Forecasting gross revenues at the movie box office," Working paper, University of Washington Seattle, WA, June 2002.
- [12]J. S. Simonoff and I. R. Sparrow, "Predicting movie grosses: Winners and losers, blockbusters and sleepers," Chance, vol. 13(3), pp. 15--24, 2000.
- [13] M. S. Sawhney and J. Eliashberg, "A parsimonious model for forecasting gross box-office revenues of motion pictures," Marketing Science, vol. Vol. 15, No. 2, pp. 113--131, 1996.
- [14] Krauss, Jonas; Nann, Stefan; Simon, Daniel; Fischbach, Kai, "Predicting movie success and academy awards through sentiment and social network analysis," University of Cologne.
- [15] P. Changkaew and R. Kongkachandra, "Automatic Movie Rating Using Visual and Linguistic Information," in NCCIT, Bangkok, Thailand, 2010.