



Movie Box Office Collection Data ETL

05.04.2019

—

DeepFlux

Three282 Technologies

Overview

One of the key projects at DeepFlux is the prediction of Box Office Collection of a live movie. To improve the model, it is important to measure accuracy of the prediction and store actual box office collection of movies. This exercise is to see the effectiveness in building code to extract, transform and load data of movie data sets. The evaluation will be based on ability to think and provision scalability in code.

Goals

1. Create Code to Crawl data from <https://boxofficecollection.in/> for various movies
2. Clean data to extract relevant information
3. Loop through multiple movies based on the config file
4. Store output in SQL table

Specifications

Box Office Collection is a site which has information of the latest Indian movies and their India Gross Collection. Their day wise prediction is available on pages with URL structure in the following format:

<https://boxofficecollection.in/luka-chuppi-box-office-collection-day-wise> which is a combination of

Base url = <https://boxofficecollection.in/>

Movie name = luka-chuppi

Tail url = [-box-office-collection-day-wise](#)

This data on the website is available in following format, where Day stands for day from release of the film

Box Office Collection



Total ₹92.05 Cr

Day 1 ₹8.01 Cr

Day 2 ₹10.08 Cr

Day 3 ₹14.04 Cr

Day 4 ₹7.90 Cr

This data has to be read for five movies (defined in a config file), and stored in a different format in a local SQL table.

Data Crawl

To crawl and parse the page HTML, BeautifulSoup library can be used. The header row of the data will have to be removed.

File Structures

Config file

The input config file will have details on movies and their release date to be used in URL generation for boxofficecollection.in and to store data in SQL table

The SQL table should have the following fields:

Movie Name

Days from Release

Date

Box Office Collection

Note: The SQL table creation should also happen via script.

3

Data Transformation

Following changes need to be made to arrive at the above columns

Movie Name- This data should store the name of movies after removing "-" from input data

Days from Release - This value should be the first field acquired from boxofficecollection.in (Day1, Day2) without any change

Date - Based on the release date, this date indicates the exact date of the collection in YYYY-MM-DD. If the movie release date is 2019-03-01 and Days from release is 2, then the date should be stored as 2019-03-02 (MM is 03, DD is 02 and YYYY is 2019 as it is the second day)

Box Office Collection - This should be stored as a float variable after taking out the currency and other signs. So 17.90 Cr will be stored as 179000000

Collating Data of Multiple Movies

The code should be written in a manner that the data of all movies stored in the config file is stored in one table with all columns as mentioned above.