# Cross-Platform Data Harmony: Machine Learning for Movie and TV ID Validation, Mapping, and Price Prediction

Rohan Shah
202218002@daiict.ac.in
MSc DS, DA-IICT
Ahmedabad, India

Internship Details:
Biztic Technologies
Ahmedabad, Gujarat
Biztic

Duration: 15 weeks
Start Date: 01st Jan 2024.
End Date: 13th April 2024

## Abstract

This project describes a research that aims to improve cross-platform data harmonization using machine learning for movie and TV ID validation, mapping, and pricing prediction. The project used proprietary code for verification, validation, and mapping, as well as machine learning methods, to simplify operations and enhance accuracy. Key components include code checking and mapping, as well as machine learning approaches for price prediction. The project's accomplishments are highlighted as efficiency benefits and accuracy improvements in data processing and prediction activities.

Keywords: Machine Learning, Verification, Validation, Mapping, Delta(Days) Prediction, SQL(Structured Query Language)

## 1 Introduction (problem statement)

The issue statement expresses the requirement to smoothly integrate and reconcile data from numerous media channels, assure proper identification and mapping of movies and TV series across these platforms, and forecast market pricing using machine learning techniques. This includes dealing with the difficulties of divergent data formats, various identification systems, and dynamic content values in the digital media world. The goal of employing machine learning algorithms is to automate these procedures, enhance accuracy, and give actionable insights for content producers, distributors, and platform owners to optimize content management, distribution strategies, and income generation.

## 2 Block Diagram (Project Structure and Flow)

**Step 1 - Identify Data Sources:** Choose the websites or online platforms from which you wish to gather information. This might include e-commerce websites, news sites, social media platforms, or any other sources that are pertinent to your research.

**Step 2 - Choose a web scraping tool:** Select an appropriate web scraping tool or library. Popular choices include BeautifulSoup (for Python), Scrapy, Selenium, and Octoparse. These tools allow you to extract data from HTML or JSON formats on websites.

**Step 3 - Understand Website Structure:** Examine the structure of the website(s) you plan to scrape. Determine which HTML elements, classes, and attributes contain the data you want. This will help you develop the best scraping approach.

**Step 4 - Create Scraping Scripts:** Write scripts to extract desired data from the recognized web pages. Navigate through the website using the web scraping tool or library of your choice, discover important items, and extract the required information. Ensure that the website's terms of service and legal requirements are followed.

**Step 5 - Data Preprocessing:** Scraped data should be cleaned and preprocessed to remove noise, handle missing values, and formatted as a structured dataset, providing good data quality for further analysis. This might involve text normalization, encoding conversion, unusual character handling, and outlier identification and removal procedures.

**Step 6 - Data Analysis:** Use exploratory data analysis (EDA) on the preprocessed dataset to get insights and analyze trends. This might include statistical summaries, data visualization with tools such as Matplotlib or Seaborn, and correlation analysis to de-

termine variables' correlations.

**Step 7 - ML Models:** A variety of machine learning models are used to acquire insights and better comprehend the information. These models assist in identifying patterns, correlations, and trends in the data. The purpose of using techniques including as classification, clustering, and dimensionality reduction is to extract useful information that may be used to inform decision-making.

**Step 8 - Regression:** Regression analysis is used to predict numerical values from input variables. In your example, a regression model is used to anticipate movie ticket prices and the number of days between a film's theatrical premiere and availability on OTT platforms. Regression algorithms allow these crucial parameters to be estimated by examining past data and taking into account aspects such as release date, genre, and market trends.

**Step 9 - Recommendation model:** Users are recommended TV series depending on their viewing history and preferences. The system uses techniques such as collaborative filtering or content-based filtering to identify similarities between persons or products (in this example, TV series) and create tailored suggestions. The cosine similarity metric is frequently used to assess the similarity of feature vectors, which aids in the selection of relevant information for each user across several platforms.
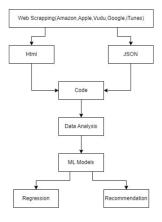


Figure 1: Flow Diagram Of Project

## 3  Tools and Technologies

### A. Visual Studio Code (VS Code)

The major tool for web scraping is Visual Studio Code (VS Code), which employs Python code written using the BeautifulSoup module. This supports data extraction from a variety of platforms, including Amazon Prime Video, Vudu, Google Play Store, and the Apple TV app.

### B. HeidiSQL

HeidiSQL is also used to connect to databases, which makes data administration and storage more efficient. Given Warner Bros. as a prominent customer, these solutions facilitate the extraction and structuring of pertinent data.

### C. Microsoft Excel 365

Microsoft Excel 365 is used to convert code outputs to.xlsx format, which ensures compatibility and simplicity of analysis. Overall, this integrated set of tools and technologies enables seamless data retrieval, management, and analysis, resulting in effective decision-making and insight development.

## 4  Use Cases

### Step 1 - Verification:

The major use case for this project is the full extraction and validation of TV show data from multiple platforms via web scraping techniques performed with BeautifulSoup. Once the data has been recovered, the following step is to thoroughly verify it for correctness. This verification step involves cross-referencing the retrieved data with information held in databases and Excel files to discover inconsistencies.

### Step 2 - Validation:

To ease this validation process, validation code is created, which allows other team members to easily identify and correct inconsistencies using a deployed internet interface. In addition, functions such as TV program merging and unmerging are provided, with checkboxes and filters used to facilitate the mapping process.



Figure 2: Validation While Opening Through API



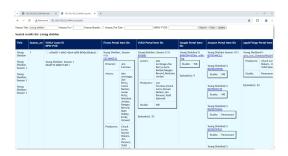Figure 3: Searching By title Representation

Figure 4: Young Sheldon Season-1.1



Figure 5: Young Sheldon Season-1.2



Figure 6: Young Sheldon Season-Bundle

**Step 3 - Mapping:**

Following that, mapping code is created to match TV program data given by Warner Bros. with platform-specific seasons, using a point system to track mapping progress and provide detailed results in Excel format. Throughout these operations, SQL databases aid data administration, guaranteeing effective processing of the project's large volume of data.
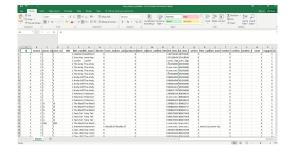


Figure 7: Season Mapping Results

## 5 Results

After methodically gathering and categorizing the dataset over three months, regression models forecast the delta days before TV series are released on various OTT platforms such as Amazon Prime Video, Vudu, iTunes, and the Apple TV app. The study seeks to predict whether TV series will be distributed one or three months following their debut. In addition, a recommendation model is created utilizing season descriptions to offer related TV series based on user feedback. Due to project challenges, experts at DAI-ICT were consulted, and it was found that clustering methods were unsuitable for the data collection. Despite the limits of some methods, the internship gave great insights into data analysis approaches and the complexity of dealing with enormous datasets. Further conversations and investigations may be necessary to improve model accuracy and efficacy in future projects.



Figure 8: Recommendation code



Figure 9: Recommendation Output

## 6 References

1. VS Code Retrieved from
   https://code.visualstudio.com/ Visual Studio

2. Web-Scrapping Retrieved from
   https://beautiful-soup-4.readthedocs.io/en/latest/
   Beautiful Soup

3. HeidiSQL. Retrieved from
   https://www.heidisql.com SQL

4. Microsoft Excel. Retrieved from
   https://www.microsoft.com/en-us/
   microsoft-365/excel

5. SQL. Retrieved from
   https://www.w3schools.com/sq1/

## 7  Github Link

Internship Project https://github.com/Rohan0412/Biztic
Biztic