



**TRƯỜNG ĐẠI HỌC BÁCH KHOA**  
**KHOA CÔNG NGHỆ THÔNG TIN**



**TIỂU LUẬN CUỐI KỲ**  
**HỌC PHẦN: KHOA HỌC DỮ LIỆU**

**ĐỀ TÀI: DỰ ĐOÁN MỨC DOANH THU PHIM**

HỌ VÀ TÊN SINH VIÊN	LỚP HỌC PHẦN	ĐIỂM BẢO VỆ
NGUYỄN TRẦN THẢO VY	20Nh91	
NGUYỄN ĐỨC QUỐC	20Nh91	
LÊ HOÀNG	20Nh91	

ĐÀ NẴNG, 05/2023

## TÓM TẮT

Dựa vào dữ liệu rất lớn của các bộ phim được tổng hợp trên website của The number, chúng em đã lựa chọn đề tài “dự đoán mức doanh thu của phim”. Sau khi nghiên cứu, nhóm đã quyết định dùng thư viện request và beautifulsoup để cào data; Label Encoder, Standard Scaling và Outliers Handling để trích xuất đặc trưng; hai mô hình là Logistic Regression và Random Forest để dự đoán doanh thu mức doanh thu phim; sau cùng là sử dụng Accuracy, Precision, Recall, F1-Score để đánh giá các mô hình. Kết quả là nhóm đã cào được Data từ web, trích xuất đặc trưng cũng như dùng hai mô hình để đánh giá mức doanh thu phim.

## BẢNG PHÂN CÔNG NHIỆM VỤ

Sinh viên thực hiện	Các nhiệm vụ	Tự đánh giá theo 3 mức (Đã hoàn thành/Chưa hoàn thành/Không triển khai)
Nguyễn Đức Quốc	<ul style="list-style-type: none"> <li>- Làm sạch dữ liệu, xử lý dữ liệu trống</li> <li>- Xử lý ngoại lệ</li> <li>- Mã hóa dữ liệu</li> <li>- Lựa chọn đặc trưng</li> </ul>	<ul style="list-style-type: none"> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> </ul>
Nguyễn Trần Thảo Vy	<ul style="list-style-type: none"> <li>- Thu thập dữ liệu</li> <li>- Thống kê mô tả trực quan về dữ liệu</li> <li>- Dán nhãn dữ liệu</li> </ul>	<ul style="list-style-type: none"> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> </ul>
Lê Hoàng	<ul style="list-style-type: none"> <li>- Mô hình hóa dữ liệu</li> <li>- Tìm bộ siêu tham số cho mô hình</li> <li>- Chuẩn hóa dữ liệu</li> <li>- Đánh giá hiệu quả mô hình</li> </ul>	<ul style="list-style-type: none"> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> <li>- Đã hoàn thành</li> </ul>

# MỤC LỤC

<b>1. GIỚI THIỆU.....</b>	<b>7</b>
1.1 CÁC VẤN ĐỀ CẦN GIẢI QUYẾT.....	7
1.2 GIẢI PHÁP.....	7
<b>2. THU THẬP VÀ MÔ TẢ DỮ LIỆU.....</b>	<b>7</b>
2.1. THU THẬP DỮ LIỆU.....	7
2.1.1. Nguồn dữ liệu.....	7
2.1.2. Công cụ thu thập.....	7
2.1.3. Cách thức thu thập.....	8
2.2. MÔ TẢ DỮ LIỆU.....	9
2.3. DÁN NHÃN DỮ LIỆU.....	12
<b>3. TRÍCH XUẤT ĐẶC TRƯNG.....</b>	<b>13</b>
3.1. LỰA CHỌN VÀ TẠO MỚI ĐẶC TRƯNG CẦN THIẾT.....	13
3.2 LÀM SẠCH DỮ LIỆU VỚI ĐÚNG KIỂU DỮ LIỆU.....	15
3.3 MÃ HÓA DỮ LIỆU.....	17
3.4 XỬ LÝ NGOẠI LỆ.....	17
3.5 CHUẨN HÓA DỮ LIỆU.....	19
3.6 ĐỘ TƯƠNG QUAN GIỮA CÁC ĐẶC TRƯNG VỚI BIẾN MỤC TIÊU.....	19
<b>4. MÔ HÌNH HÓA DỮ LIỆU.....</b>	<b>20</b>
4.1 CÁC METRICS DÙNG ĐỂ ĐÁNH GIÁ:.....	20
4.2 MÔ HÌNH LOGISTIC REGRESSION:.....	21
4.3 MÔ HÌNH RANDOM FOREST.....	24
<b>5. KẾT LUẬN.....</b>	<b>26</b>
<b>6. TÀI LIỆU THAM KHẢO.....</b>	<b>26</b>

# DANH MỤC HÌNH ẢNH

Hình 1. Thẻ chứa doanh thu toàn cầu của một bộ phim.....	8
Hình 2. Thẻ chứa link đến thông tin chi tiết của một bộ phim.....	8
Hình 3. Phân tích cấu trúc HTML của một bộ phim.....	9
Hình 4. Dữ liệu BigDS_Raw thu được.....	11
Hình 5. Dữ liệu SmallDS_clean thu được sau khi clean dữ liệu.....	12
Hình 6. Dữ liệu doanh thu phim trước và sau khi dán nhãn dữ liệu.....	13
Hình 7. Trích xuất các đặc trưng lựa chọn trong web.....	13
Hình 8. Trích xuất các đặc trưng lựa chọn trong web.....	14
Hình 9. Các đặc trưng WorldwideBox Office, Production Budget, Date Release sau khi xử lý .....	16
Hình 10. Các đặc trưng WorldwideBox Office, Production Budget sau khi xử lý kiểu dữ liệu .....	16
Hình 11. Dữ liệu huấn luyện trước khi được mã hóa.....	17
Hình 12. Dữ liệu huấn luyện sau khi được mã hóa.....	17
Hình 13. Phân bố dữ liệu huấn luyện của đặc trưng Production Budget khi chưa xử lý ngoại lệ.....	18
Hình 14. Phân bố dữ liệu huấn luyện của đặc trưng Running Time khi chưa xử lý ngoại lệ..	18
Hình 15. Phân bố dữ liệu huấn luyện sau khi xử lý ngoại lệ.....	19
Hình 16. Dữ liệu huấn luyện sau khi chuẩn hóa dữ liệu.....	19
Hình 17. Sự tương quan giữa các đặc trưng với biến mục tiêu.....	19
Hình 18. Sự tương quan giữa các đặc trưng với biến mục tiêu.....	20
Hình 19. Độ tương quan giảm dần của các đặc trưng so với biến mục tiêu.....	20
Hình 20. Đồ thị thể hiện hiệu suất 4 metrics trên mô hình logistic regresstion.....	23
Hình 21. Ma trận nhầm lẫn kết quả dự đoán với kết quả thực của test của mô hình Logistic Regresstion.....	23
Hình 22. Đánh giá mô hình Logistic Regresstion theo các metrics là Accuary, Precesion, Recall, F1-Scall.....	24

Hình 23. Ma trận nhầm lẫn kết quả dự đoán với kết quả thực của test của mô hình Random Forest.....	25
Hình 24. Đồ thị thể hiện hiệu suất của 4 metrics trên mô hình Random Forest.....	25
Hình 25. Đánh giá mô hình theo các metrics là Accuracy, Precision, Recall, F1-Score theo Random Forest.....	26

# 1. Giới thiệu

## 1.1 Các vấn đề cần giải quyết

- Làm sao để cào được dữ liệu từ web.
- Cách để trích xuất dữ liệu
- Phương pháp để dự đoán doanh thu phim và độ chính xác của phương pháp đó

## 1.2 Giải pháp

- Sử dụng thư viện request, beautifulsoup phục vụ cho việc cào dữ liệu.
- Quan sát bằng mắt và chọn lọc, sau đó cào các đặc trưng đã chọn tương ứng với các thẻ HTML.
- Sử dụng hai mô hình học máy là RandomForest và Linear Regression.

# 2. Thu thập và mô tả dữ liệu

## 2.1. Thu thập dữ liệu

### 2.1.1. Nguồn dữ liệu

Dữ liệu được thu thập từ trang web thống kê, đánh giá về ngành công nghiệp điện ảnh và truyền hình thế giới: <https://www.the-numbers.com/>

### 2.1.2. Công cụ thu thập

- IDE: Visual Studio Code
- Ngôn ngữ: Python
- Thư viện: beautifulsoup4, request

## 2.1.3. Cách thức thu thập

### Bước 1: Lấy tất cả doanh thu và link thông tin chi tiết của các bộ phim

Other Worldwide Cumulative records: All Time Single Market - All Time Animated Worldwide - All Time Sequel Worldwide - All Time Non-Sequel Worldwide - Top 2022 Worldwide - Top 2023 Worldwide

This chart contains the top movies based on the cumulative worldwide box office.

Rank	Year	Movie	Worldwide Box Office	Domestic Box Office	International Box Office
1	2009	Avatar	\$2,923,706,026	\$765,221,649	\$2,138,484,377
2	2019	Avengers: Endgame	\$2,798,371,755	\$658,373,000	\$1,936,358,755
3	2022	Avatar: The Way of Water	\$2,319,740,227	\$684,052,555	\$1,635,687,672
4	1997	Titanic	\$2,201,958,563	\$584,360,795	\$1,548,588,773
5	2015	Star Wars Ep. VII: The Force Awakens	\$2,069,615,917	\$936,662,225	\$1,127,953,592
6	2018	Avengers: Infinity War	\$2,048,359,754	\$678,815,482	\$1,369,544,272
7	2021	Spider-Man: No Way Home	\$1,910,048,245	\$814,115,070	\$1,095,933,175
8	2015	Jurassic World	\$1,669,593,641	\$652,306,625	\$1,017,657,016
9	2019	The Lion King	\$1,647,733,638	\$543,638,043	\$1,104,095,595
10	2012	The Avengers	\$1,515,100,211	\$623,357,910	\$891,742,301
11	2015	Furious 7	\$1,511,886,364	\$353,007,020	\$1,158,979,344
12	2022	Top Gun: Maverick	\$1,481,230,730	\$718,732,821	\$762,497,909
13	2019	Frozen II	\$1,453,683,476	\$477,373,578	\$976,309,898

Additional DevTools components have been installed. Reload for full DevTools capabilities. [Reload DevTools](#)

Elements

```
<table>
  <tr>
    <td class="data">1</td>
    <td class="data">2009</td>
    <td class="data">Avatar</td>
    <td align="right">$2,923,706,026</td>
    <td align="right">$765,221,649</td>
    <td align="right">$2,138,484,377</td>
  </tr>
  <tr>
    <td class="highlight">2</td>
    <td class="highlight">2019</td>
    <td class="highlight">Avengers: Endgame</td>
    <td align="right">$2,798,371,755</td>
    <td align="right">$658,373,000</td>
    <td align="right">$1,936,358,755</td>
  </tr>
  <tr>
    <td class="highlight">3</td>
    <td class="highlight">2022</td>
    <td class="highlight">Avatar: The Way of Water</td>
    <td align="right">$2,319,740,227</td>
    <td align="right">$684,052,555</td>
    <td align="right">$1,635,687,672</td>
  </tr>
  <tr>
    <td class="highlight">4</td>
    <td class="highlight">1997</td>
    <td class="highlight">Titanic</td>
    <td align="right">$2,201,958,563</td>
    <td align="right">$584,360,795</td>
    <td align="right">$1,548,588,773</td>
  </tr>
  <tr>
    <td class="highlight">5</td>
    <td class="highlight">2015</td>
    <td class="highlight">Star Wars Ep. VII: The Force Awakens</td>
    <td align="right">$2,069,615,917</td>
    <td align="right">$936,662,225</td>
    <td align="right">$1,127,953,592</td>
  </tr>
  <tr>
    <td class="highlight">6</td>
    <td class="highlight">2018</td>
    <td class="highlight">Avengers: Infinity War</td>
    <td align="right">$2,048,359,754</td>
    <td align="right">$678,815,482</td>
    <td align="right">$1,369,544,272</td>
  </tr>
  <tr>
    <td class="highlight">7</td>
    <td class="highlight">2021</td>
    <td class="highlight">Spider-Man: No Way Home</td>
    <td align="right">$1,910,048,245</td>
    <td align="right">$814,115,070</td>
    <td align="right">$1,095,933,175</td>
  </tr>
  <tr>
    <td class="highlight">8</td>
    <td class="highlight">2015</td>
    <td class="highlight">Jurassic World</td>
    <td align="right">$1,669,593,641</td>
    <td align="right">$652,306,625</td>
    <td align="right">$1,017,657,016</td>
  </tr>
  <tr>
    <td class="highlight">9</td>
    <td class="highlight">2019</td>
    <td class="highlight">The Lion King</td>
    <td align="right">$1,647,733,638</td>
    <td align="right">$543,638,043</td>
    <td align="right">$1,104,095,595</td>
  </tr>
  <tr>
    <td class="highlight">10</td>
    <td class="highlight">2012</td>
    <td class="highlight">The Avengers</td>
    <td align="right">$1,515,100,211</td>
    <td align="right">$623,357,910</td>
    <td align="right">$891,742,301</td>
  </tr>
  <tr>
    <td class="highlight">11</td>
    <td class="highlight">2015</td>
    <td class="highlight">Furious 7</td>
    <td align="right">$1,511,886,364</td>
    <td align="right">$353,007,020</td>
    <td align="right">$1,158,979,344</td>
  </tr>
  <tr>
    <td class="highlight">12</td>
    <td class="highlight">2022</td>
    <td class="highlight">Top Gun: Maverick</td>
    <td align="right">$1,481,230,730</td>
    <td align="right">$718,732,821</td>
    <td align="right">$762,497,909</td>
  </tr>
  <tr>
    <td class="highlight">13</td>
    <td class="highlight">2019</td>
    <td class="highlight">Frozen II</td>
    <td align="right">$1,453,683,476</td>
    <td align="right">$477,373,578</td>
    <td align="right">$976,309,898</td>
  </tr>
</table>
```

Annual Box Office  
Box Office Records  
International Box Office  
Distributors  
People Records  
People Index  
Genre Tracking  
Keyword Tracking  
Franchises  
Research Tools  
Bankability Index

Most Anticipated Movies

The Little Mermaid  
Indiana Jones and the Dial of Destiny  
The Machine  
Mission: Impossible Dead Reckoning Part One  
Spider-Man: Across the Spider-Verse  
Transformers: Rise of the Beasts  
Five Nights at Freddy's  
About My Father  
Kandahar  
The Flash

Hình 1. Thẻ chứa doanh thu toàn cầu của một bộ phim

Other Worldwide Cumulative records: All Time Single Market - All Time Animated Worldwide - All Time Sequel Worldwide - All Time Non-Sequel Worldwide - Top 2022 Worldwide - Top 2023 Worldwide

This chart contains the top movies based on the cumulative worldwide box office.

Rank	Year	Movie	Worldwide Box Office	Domestic Box Office	International Box Office
1	2009	Avatar	\$2,923,706,026	\$765,221,649	\$2,138,484,377
2	2019	Avengers: Endgame	\$2,798,371,755	\$658,373,000	\$1,936,358,755
3	2022	Avatar: The Way of Water	\$2,319,740,227	\$684,052,555	\$1,635,687,672
4	1997	Titanic	\$2,201,958,563	\$584,360,795	\$1,548,588,773
5	2015	Star Wars Ep. VII: The Force Awakens	\$2,069,615,917	\$936,662,225	\$1,127,953,592
6	2018	Avengers: Infinity War	\$2,048,359,754	\$678,815,482	\$1,369,544,272
7	2021	Spider-Man: No Way Home	\$1,910,048,245	\$814,115,070	\$1,095,933,175
8	2015	Jurassic World	\$1,669,593,641	\$652,306,625	\$1,017,657,016
9	2019	The Lion King	\$1,647,733,638	\$543,638,043	\$1,104,095,595
10	2012	The Avengers	\$1,515,100,211	\$623,357,910	\$891,742,301
11	2015	Furious 7	\$1,511,886,364	\$353,007,020	\$1,158,979,344
12	2022	Top Gun: Maverick	\$1,481,230,730	\$718,732,821	\$762,497,909
13	2019	Frozen II	\$1,453,683,476	\$477,373,578	\$976,309,898

Additional DevTools components have been installed. Reload for full DevTools capabilities. [Reload DevTools](#)

Elements

```
<table>
  <tr>
    <td class="data">1</td>
    <td class="data">2009</td>
    <td class="data">Avatar</td>
    <td align="right">$2,923,706,026</td>
    <td align="right">$765,221,649</td>
    <td align="right">$2,138,484,377</td>
  </tr>
  <tr>
    <td class="data">2</td>
    <td class="data">2019</td>
    <td class="data">Avengers: Endgame</td>
    <td align="right">$2,798,371,755</td>
    <td align="right">$658,373,000</td>
    <td align="right">$1,936,358,755</td>
  </tr>
  <tr>
    <td class="data">3</td>
    <td class="data">2022</td>
    <td class="data">Avatar: The Way of Water</td>
    <td align="right">$2,319,740,227</td>
    <td align="right">$684,052,555</td>
    <td align="right">$1,635,687,672</td>
  </tr>
  <tr>
    <td class="data">4</td>
    <td class="data">1997</td>
    <td class="data">Titanic</td>
    <td align="right">$2,201,958,563</td>
    <td align="right">$584,360,795</td>
    <td align="right">$1,548,588,773</td>
  </tr>
  <tr>
    <td class="data">5</td>
    <td class="data">2015</td>
    <td class="data">Star Wars Ep. VII: The Force Awakens</td>
    <td align="right">$2,069,615,917</td>
    <td align="right">$936,662,225</td>
    <td align="right">$1,127,953,592</td>
  </tr>
  <tr>
    <td class="data">6</td>
    <td class="data">2018</td>
    <td class="data">Avengers: Infinity War</td>
    <td align="right">$2,048,359,754</td>
    <td align="right">$678,815,482</td>
    <td align="right">$1,369,544,272</td>
  </tr>
  <tr>
    <td class="data">7</td>
    <td class="data">2021</td>
    <td class="data">Spider-Man: No Way Home</td>
    <td align="right">$1,910,048,245</td>
    <td align="right">$814,115,070</td>
    <td align="right">$1,095,933,175</td>
  </tr>
  <tr>
    <td class="data">8</td>
    <td class="data">2015</td>
    <td class="data">Jurassic World</td>
    <td align="right">$1,669,593,641</td>
    <td align="right">$652,306,625</td>
    <td align="right">$1,017,657,016</td>
  </tr>
  <tr>
    <td class="data">9</td>
    <td class="data">2019</td>
    <td class="data">The Lion King</td>
    <td align="right">$1,647,733,638</td>
    <td align="right">$543,638,043</td>
    <td align="right">$1,104,095,595</td>
  </tr>
  <tr>
    <td class="data">10</td>
    <td class="data">2012</td>
    <td class="data">The Avengers</td>
    <td align="right">$1,515,100,211</td>
    <td align="right">$623,357,910</td>
    <td align="right">$891,742,301</td>
  </tr>
  <tr>
    <td class="data">11</td>
    <td class="data">2015</td>
    <td class="data">Furious 7</td>
    <td align="right">$1,511,886,364</td>
    <td align="right">$353,007,020</td>
    <td align="right">$1,158,979,344</td>
  </tr>
  <tr>
    <td class="data">12</td>
    <td class="data">2022</td>
    <td class="data">Top Gun: Maverick</td>
    <td align="right">$1,481,230,730</td>
    <td align="right">$718,732,821</td>
    <td align="right">$762,497,909</td>
  </tr>
  <tr>
    <td class="data">13</td>
    <td class="data">2019</td>
    <td class="data">Frozen II</td>
    <td align="right">$1,453,683,476</td>
    <td align="right">$477,373,578</td>
    <td align="right">$976,309,898</td>
  </tr>
</table>
```

Annual Box Office  
Box Office Records  
International Box Office  
Distributors  
People Records  
People Index  
Genre Tracking  
Keyword Tracking  
Franchises  
Research Tools  
Bankability Index

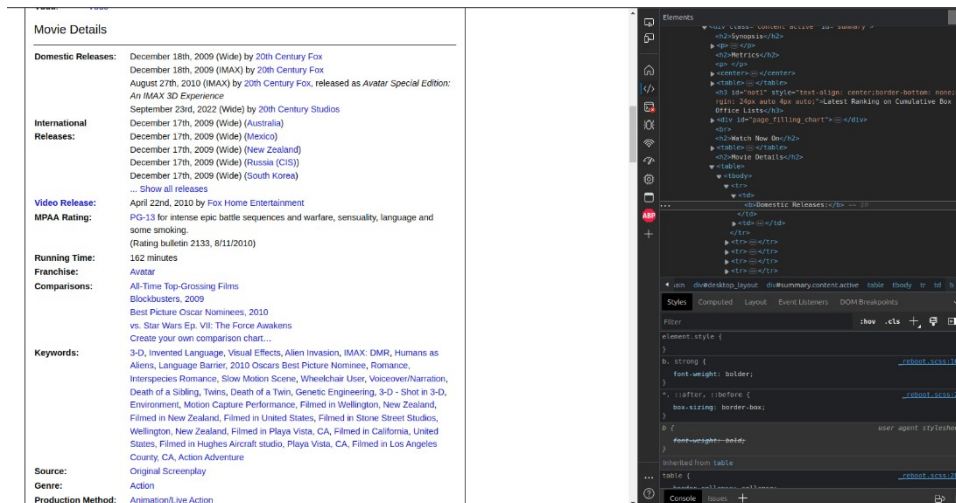
Most Anticipated Movies

The Little Mermaid  
Indiana Jones and the Dial of Destiny  
The Machine  
Mission: Impossible Dead Reckoning Part One  
Spider-Man: Across the Spider-Verse  
Transformers: Rise of the Beasts  
Five Nights at Freddy's  
About My Father  
Kandahar  
The Flash

Hình 2. Thẻ chứa link đến thông tin chi tiết của một bộ phim



## Bước 2: Phân tích cấu trúc trang web của từng bộ phim



Hình 3. Phân tích cấu trúc HTML của một bộ phim

- Tìm kiếm các thẻ chứa thông tin cần thiết của bộ phim.

### Bước 3: Trích xuất thông tin

- Sử dụng thư viện request để get trang web.
- Đọc HTML của trang web thông qua hàm của thư viện beautifulsoup.

## Bước 4: Chuyển list dữ liệu sang file csv

- Input : Link web.
- Output : 1 file raw csv.

## 2.2. Mô tả dữ liệu

- Số mẫu :
  - SmallDS : 1000 (samples).
  - BigDS : 10000 (samples).
- Chiều dữ liệu : 10000 x 16 (gồm 15 đặc trưng và một cột số thứ tự).
  - SmallDS : 1000 x 16 (gồm 15 đặc trưng và một cột số thứ tự).
  - BigDS : 10000 x 16 (gồm 15 đặc trưng và một cột số thứ tự).
- Số đặc trưng của mẫu: 15 đặc trưng.

Bảng 1: Mô tả dữ liệu thô ban đầu của BigDS

Column	Kiểu dữ liệu	Số dữ liệu trống
Rank	Int64	0
Year	Int64	0
Movie	object	0
WorldwideBox Office	object	0
Production Budget	object	5293
Date Releases	object	4710
MPAA	object	2510
Running Time	object	1431
Franchise	object	8036
Genre	object	489
Creative Type	object	739
Production/Financing Companies	object	5329
Production Countries	object	443
Languages	object	1763
Director	object	1297

Bảng 2: Mô tả dữ liệu thô ban đầu của SmallDS

Column	Kiểu dữ liệu	Số dữ liệu trống
Rank	Int64	0
Year	Int64	0
Movie	object	0
WorldwideBox Office	object	0
Production Budget	object	83
Date Releases	object	255
MPAA	object	28
Running Time	object	12
Franchise	object	429
Genre	object	1
Creative Type	object	1
Production/Financing Companies	object	106
Production Countries	object	2
Languages	object	13
Director	object	15

- Dữ liệu SmallDS\_raw thu thập được:

7df\_small.head()

✓0.0s

Hình 4. Dữ liệu BigDS\_Raw thu được

- Dữ liệu SmallDS\_clean thu được sau khi clean dữ liệu

Unnamed: 0	Rank	Year	Movie	WorldwideBox Office	Production Budget	Date Releases	MPAA	Running Time	Franchise	Genre	Creative Type	Production/Financing Companies	Pre C
0	0	1	2009	Avatar	2923706026	237000000.0	December 17th, 2009	PG-13	162.0	Avatar	Action	Science Fiction	Dune Entertainment, 20th Century Fox, Ingeniou...
1	1	2	2019	Avengers-Endgame-(2019)	2794731755	400000000.0	April 23rd, 2019	PG-13	181.0	Marvel Cinematic UniverseAvengers	Action	Super Hero	Marvel Studios
2	2	3	2022	Avatar-The-Way-of-Water-(2022)	2319738066	460000000.0	December 9th, 2022	PG-13	190.0	Avatar	Action	Science Fiction	Lightstorm Entertainment, 20th Century Studios...
3	3	4	1997	Titanic-(1997)	2222985568	200000000.0	December 18th, 1997	PG-13	194.0	NaN	Drama	Historical Fiction	20th Century Fox, Paramount Pictures, Lightsto...
4	4	5	2015	Star-Wars-Ep-VII-The-Force-Awakens	2064615817	306000000.0	December 16th, 2015	PG-13	136.0	Star Wars	Adventure	Science Fiction	Lucasfilm, Bad Robot
...	...	...	...	...	...	...	...	...	...	...	...	...	...
995	995	996	1995	Dangerous-Minds	178919401	23000000.0	NaN	R	99.0	NaN	Drama	Dramatization	NaN
996	996	997	2006	Scary-Movie-4-(2006)	178710620	40000000.0	NaN	PG-13	83.0	Scary Movie	Comedy	Contemporary Fiction	NaN
997	997	998	2003	How-to-Lose-a-Guy-in-10-	178503788	50000000.0	April 18th, 2003	PG-13	116.0	NaN	Romantic Comedy	Contemporary Fiction	NaN

Hình 5. Dữ liệu SmallDS\_clean thu được sau khi clean dữ liệu

## 2.3. Dán nhãn dữ liệu

- Doanh thu phim sẽ được dán nhãn với các mức doanh thu phim tương ứng.

Mức doanh thu (triệu đô)	Nhãn tương ứng
0-300	1
300-600	2
600-900	3
900-1200	4
1200-1500	5
1500-1800	6
1800-2100	7
2100-2400	8
2400-2700	9
2700-3000	10

Bảng 3: Bảng đánh giá mức doanh thu phim

WorldwideBox Office		WorldwideBox Office	
105	781947691	105	3
68	894860230	68	3
479	322459006	479	2
399	362605033	399	2
434	348901032	434	2
...	...	...	...
835	209221328	835	1
192	560483719	192	2
629	258751370	629	1
559	288510892	559	1
684	245179530	684	1

Hình 6. Dữ liệu doanh thu phim trước và sau khi dán nhãn dữ liệu

### 3. Trích xuất đặc trưng

#### 3.1. Lựa chọn và tạo mới đặc trưng cần thiết

**Theatrical Performance**

Domestic Box Office	\$785,221,649	<a href="#">Details</a>
International Box Office	\$2,138,484,377	<a href="#">Details</a>
Worldwide Box Office	\$2,923,706,026	

**Home Market Performance**

Est. Domestic DVD Sales	\$209,780,313	<a href="#">Details</a>
Est. Domestic Blu-ray Sales	\$220,499,602	<a href="#">Details</a>
Total Est. Domestic Video Sales	\$430,279,915	

[Further financial details...](#)

[Summary](#)
[News](#)
[Box Office](#)
[International](#)
[Video Sales](#)
[Full Financials](#)
[Cast & Crew](#)
[Trailer](#)

**Synopsis**

Jake Sully is a wounded ex-marine, thrust into an effort to settle and exploit Pandora, an exotic moon rich in bio-diversity and inhabited by the Na'vi, a ten-foot-tall humanoid species. After Neytiri, a female Na'vi, rescues Jake after he becomes separated from his team, he learns more about the planet and eventually crosses over to lead the indigenous race in a battle for survival.

**Metrics**

**Opening Weekend:** \$77,025,481 (9.8% of total gross)

**Legs:** 10.19 (domestic box office/biggest weekend)

**Domestic Share:** 26.9% (domestic box office/worldwide)

**Production Budget:** \$237,000,000 (worldwide box office is 12.3 times production budget)

**Theater counts:** 3,452 opening theaters/3,461 max. theaters, 14.8 weeks average run per theater

**Infl. Adj. Dom. BO** \$929,879,724

**Quick Links**

- DEG Watched at Home Top 20
- Netflix Daily Top 10
- Weekly DVD+Blu-ray Chart
- News
- Release Schedule
- Daily Box Office
- Weekend Box Office
- Weekly Box Office
- Annual Box Office
- Box Office Records
- International Box Office
- Distributors
- People Records
- People Index
- Genre Tracking
- Keyword Tracking
- Franchises
- Research Tools
- Bankability Index

**Most Anticipated Movies**

- The Little Mermaid
- Come Out Fighting
- Mission: Impossible Dead Reckoning Part One
- The Machine
- About My Father
- Five Nights at Freddy's
- Spider-Man: Across the Spider-Verse
- Transformers: Rise of the Beasts
- The Flash
- Kandahar

**Trending Movies**

- The Super Mario Bros. Movie
- Guardians of the Galaxy Vol 3
- John Wick: Chapter 4
- Dungeons & Dragons: Honor Among Thieves
- Avatar: The Way of Water
- Evil Dead Rise
- Ant-Man and the Wasp: Quantumania
- Fast X
- Love Again
- Book Club: The Next Chapter

Hình 7. Trích xuất các đặc trưng lựa chọn trong web

Watch Now On		Arnold Schwarzen...
<div><div><div><div><div></div><div>iTunes:</div><div>iTunes, iTunes</div></div></div><div><div><div><div></div><div>Google Play:</div><div>Google Play, Google Play, Google Play</div></div></div><div><div><div><div></div><div>Vudu:</div><div>Vudu</div></div></div></div></div></div></div>		
Movie Details		
<div><div><div><div><div>Domestic Releases:</div><div>December 18th, 2009 (Wide) by 20th Century Fox</div><div>December 18th, 2009 (IMAX) by 20th Century Fox</div><div>August 27th, 2010 (IMAX) by 20th Century Fox, released as <i>Avatar Special Edition: An IMAX 3D Experience</i></div><div>September 23rd, 2022 (Wide) by 20th Century Studios</div></div></div><div><div><div><div>International Releases:</div><div>December 17th, 2009 (Wide) (Australia)</div><div>December 17th, 2009 (Wide) (Mexico)</div><div>December 17th, 2009 (Wide) (New Zealand)</div><div>December 17th, 2009 (Wide) (Russia (CIS))</div><div>December 17th, 2009 (Wide) (South Korea)</div><div>... Show all releases</div></div></div><div><div><div><div>Video Release:</div><div>April 22nd, 2010 by Fox Home Entertainment</div></div></div><div><div><div><div>MPAA Rating:</div><div>PG-13 for intense epic battle sequences and warfare, sensuality, language and some smoking.</div><div>(Rating bulletin 2133, 8/11/2010)</div></div></div><div><div><div><div>Running Time:</div><div>162 minutes</div></div></div><div><div><div><div>Franchise:</div><div>Avatar</div></div></div><div><div><div><div>Comparisons:</div><div>All-Time Top-Grossing Films</div><div>Blockbusters, 2009</div><div>Best Picture Oscar Nominees, 2010</div><div>vs. Star Wars Ep. VII: The Force Awakens</div><div>Create your own comparison chart...</div></div></div><div><div><div><div>Keywords:</div><div>3-D, Invented Language, Visual Effects, Alien Invasion, IMAX: DMR, Humans as Aliens, Language Barrier, 2010 Oscars Best Picture Nominee, Romance, Interspecies Romance, Slow Motion Scene, Wheelchair User, Voiceover/Narration, Death of a Sibling, Twins, Death of a Twin, Genetic Engineering, 3-D - Shot in 3-D, Environment, Motion Capture Performance, Filmed in Wellington, New Zealand, Filmed in New Zealand, Filmed in United States, Filmed in Stone Street Studios, Wellington, New Zealand, Filmed in Playa Vista, CA, Filmed in California, United States, Filmed in Hughes Aircraft studio, Playa Vista, CA, Filmed in Los Angeles County, CA, Action Adventure</div></div></div><div><div><div><div>Source:</div><div>Original Screenplay</div></div></div><div><div><div><div>Genre:</div><div>Action</div></div></div><div><div><div><div>Production Method:</div><div>Animation/Live Action</div></div></div><div><div><div><div>Creative Type:</div><div>Science Fiction</div></div></div><div><div><div><div>Production/Financing</div><div>Dune Entertainment, 20th Century Fox, Ingenious Film Partners</div></div></div><div><div><div><div>Companies:</div></div></div><div><div><div><div>Production Countries:</div><div>United States</div></div></div><div><div><div><div>Languages:</div><div>English, Na'vi</div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>		
<div><div><div><div><div>Production and Technical Credits</div><div><div><div><div>James Cameron</div><div>Director</div></div><div><div><div>James Cameron</div><div>Screenwriter</div></div><div><div><div>James Cameron</div><div>Producer</div></div><div><div><div>Jon Landau</div><div>Producer</div></div><div><div><div>Colin Wilson</div><div>Executive Producer</div></div><div><div><div>Laeta Kalogridis</div><div>Executive Producer</div></div><div><div><div>Brooke Breton</div><div>Co-Producer</div></div><div><div><div>Josh McLaglen</div><div>Co-Producer</div></div></div></div><div><div><div><div>Stephen Rivkin</div><div>Editor</div></div><div><div><div>Robert Stromberg</div><div>Production Designer</div></div><div><div><div>Rick Carter</div><div>Production Designer</div></div><div><div><div>Mauro Fiore</div><div>Cinematographer</div></div><div><div><div>John Refoua</div><div>Editor</div></div><div><div><div>James Horner</div><div>Composer</div></div><div><div><div>Yuri Bartoli</div><div>Supervising Visual Art Director</div></div><div><div><div>Kim Sinclair</div><div>Lead Supervising Art Director</div></div></div></div><div><div><div><div>Kevin Ishioka</div><div>Supervising Art Director</div></div><div><div><div>Stefan Dechant</div><div>Supervising Art Director</div></div><div><div><div>Todd Cherniawsky</div><div>Supervising Art Director</div></div><div><div><div>Andrew L. Jones</div><div>Virtual Production Art Director</div></div><div><div><div>Norm Newberry</div><div>Art Director</div></div><div><div><div>Nick Bassett</div><div>Art Director</div></div><div><div><div>Rob Bavin</div><div>Art Director</div></div><div><div><div>Simon Bright</div><div>Art Director</div></div><div><div><div>Jill Cormack</div><div>Art Director</div></div><div><div><div>Sean Haworth</div><div>Art Director</div></div><div><div><div>Mayes C. Rubeo</div><div>Costume Designer</div></div><div><div><div>Andrew Menzies</div><div>Art Director</div></div><div><div><div>Deborah L. Scott</div><div>Costume Designer</div></div><div><div><div>Andy McLaren</div><div>Art Director</div></div><div><div><div>Jim Tanenbaum</div><div>Sound Mixer</div></div><div><div><div>Christopher Boyes</div><div>Supervising Sound Editor</div></div><div><div><div>Christopher Boyes</div><div>Sound Designer</div></div><div><div><div>Shannon Mills</div><div>Sound Effects Editor</div></div><div><div><div>Christopher Boyes</div><div>Re-recording Mixer</div></div><div><div><div>Gary Summers</div><div>Re-recording Mixer</div></div><div><div><div>Joe Letteri</div><div>Visual Effects Supervisor</div></div><div><div><div>Andy Nelson</div><div>Re-recording Mixer</div></div><div><div><div>Stephen Rosenbaum</div><div>Weta visual effects supervisor</div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>		

Hình 8. Trích xuất các đặc trưng lựa chọn trong web

Các đặc trưng được lựa chọn: Rank, Year, Movie, WorldwideBox Office, Production Budget, Date Releases, MPAA, Running Time, Franchise, Genre, Creative Type, Production/Financing Companies, Production Countries, Languages, Director

Các đặc trưng không được lựa chọn: Reviews, Score, Rating, Population. Vì mục tiêu của nhóm là dự đoán mức doanh thu khi phim còn trong giai đoạn ý tưởng, hay vừa mới ra mắt, các đặc trưng này chưa tồn tại ở thời điểm đó.

Các đặc trưng được tạo thêm: Từ Date Release tạo ra hai đặc trưng mới là Day Release và Month Release.

### **3.2 Làm sạch dữ liệu với đúng kiểu dữ liệu**

- Sau khi đã lựa chọn các đặc trưng thì cào dữ liệu theo các thẻ HTML tương ứng.
- Loại bỏ các giá trị NaN bằng hàm mode().
- Loại bỏ các kí tự đặc biệt và chữ trong cột: WorldwideBox Office, Production Budget, Date Releases, Running Time.
- Các đặc trưng WorldwideBox Office, Production Budget sau khi xử lý và từ đặc trưng Date Releases tạo thêm hai đặc trưng mới Day Releases và Month Release

Rank	Year	Movie	WorldwideBox Office	Production Budget	Date Releases	MPAA	Running Time	Franchise	Day Releases	Month Releases
1	2009	Avatar	2923706026	2370000000.0	December 17th, 2009	PG-13	162.0	5	3.0	12.0
2	2019	Avengers-Endgame-(2019)	2794731755	4000000000.0	April 23rd, 2019	PG-13	181.0	12	1.0	4.0
3	2022	Avatar-The-Way-of-Water-(2022)	2319738066	4600000000.0	December 9th, 2022	PG-13	190.0	5	4.0	12.0
4	1997	Titanic-(1997)	2222985568	2000000000.0	December 18th, 1997	PG-13	194.0	5	3.0	12.0
5	2015	Star-Wars-Ep-VII-The-Force-Awakens	2064615817	3060000000.0	December 16th, 2015	PG-13	136.0	5	2.0	12.0

Hình 9. Các đặc trưng WorldwideBox Office, Production Budget, Date Release sau khi xử lý

- Các đặc trưng Rank, WorldwideBox Office, Production Budget, Running Time sau khi xử lý kiểu dữ liệu

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 18 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   Unnamed: 0                            1000 non-null   int64  
 1   Rank                                  1000 non-null   object  
 2   Year                                  1000 non-null   int64  
 3   Movie                                 1000 non-null   object  
 4   WorldwideBox Office                    1000 non-null   int64  
 5   Production Budget                      917 non-null    float64 
 6   Date Releases                         745 non-null    object  
 7   MPAA                                  972 non-null    object  
 8   Running Time                          988 non-null    float64 
 9   Franchise                             571 non-null    object  
10   Genre                                 999 non-null    object  
11   Creative Type                         999 non-null    object  
12   Production/Financing Companies         894 non-null    object  
13   Production Countries                  998 non-null    object  
14   Languages                             987 non-null    object  
15   Director                             985 non-null    object  
16   Day Releases                          745 non-null    float64 
17   Month Releases                       745 non-null    float64 
dtypes: float64(4), int64(3), object(11)
memory usage: 140.8+ KB
```

Hình 10. Các đặc trưng WorldwideBox Office, Production Budget sau khi xử lý kiểu dữ liệu



### 3.3 Mã hóa dữ liệu

- Các bộ phim có đạo diễn thuộc top 100 đạo diễn có doanh thu cao nhất sẽ được gán là 1 ngược lại là 0.
- Các bộ phim có công ty sản xuất thuộc top 100 công ty có doanh thu cao nhất sẽ được gán là 1 ngược lại là 0.
- Các bộ phim thuộc top 15 series có doanh thu cao nhất sẽ được gán là 1 ngược lại là 0, nếu không thuộc series nào sẽ được gán là 2.
- Các đặc trưng kiểu category khác được mã hóa bằng LabelEncoder.

Rank	Year	Movie	WorldwideBox Office	Production Budget	Date Releases	MPAA	Running Time	Franchise	Genre	Creative Type	Production/Financing Companies	Production Countries	Languages	Director	Day Releases	Month Releases
106	2016	Deadpool	781947691	58000000.0	February 9th, 2016	R	107.0	X-MenDeadpool	Action	Super Hero	[Marvel Studios]	[United States]	[English]	[Tim Miller]	1.0	2.0
69	2007	Spider-Man-3	894860230	258000000.0	NaN	PG-13	139.0	Spider-Man	Adventure	Super Hero	[Columbia Pictures, Marvel Studios, Laura Zisk...	[United States]	[English]	[Sam Raimi]	NaN	NaN
480	2010	Robin-Hood-(2010)	322459006	210000000.0	NaN	PG-13	139.0	NaN	Action	Historical Fiction	[United Artists, Fairbanks]	[United Kingdom, United States]	[English]	[Ridley Scott]	NaN	NaN
400	2007	Alvin-and-the-Chipmunks-(2007)	362605033	55000000.0	NaN	PG	91.0	Alvin and the Chipmunks	Adventure	Kids Fiction	[Regency Enterprises, Fox 2000 Pictures, Bagda...	[United States]	[English]	[Tim Hill]	NaN	NaN
435	2018	Mary-Poppins-Returns-(2018)	348901032	130000000.0	December 20th, 2018	PG	130.0	Mary Poppins	Musical	Kids Fiction	[Walt Disney Pictures, Lucamar, Marc Platt Pro...	[United States]	[English]	[Rob Marshall]	3.0	12.0

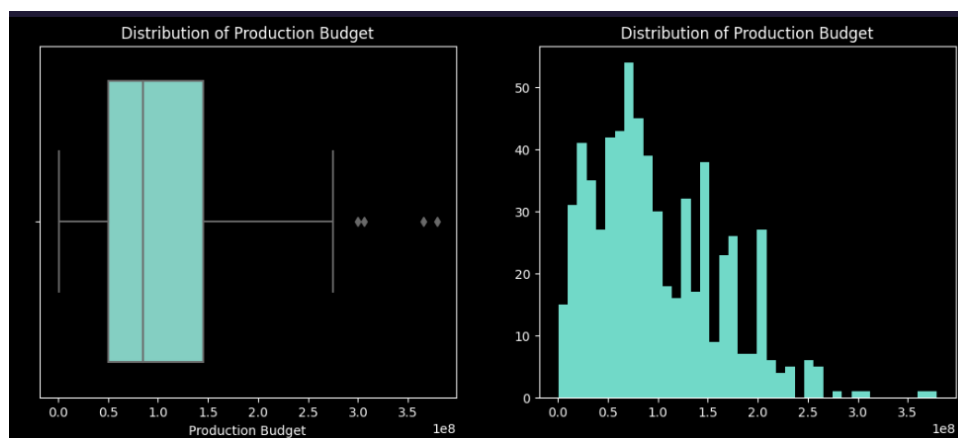
Hình 11. Dữ liệu huấn luyện trước khi được mã hóa

	Production Budget	MPAA	Running Time	Genre	Creative Type	topCompanies	topDirector	Day Releases	Month Releases	holiday_month	origin_isEnglish	topFranchise	countries_isUS
105	58000000.0	7	107.0	0	8	1	1	1.0	2.0	0	1	2	1
68	258000000.0	4	139.0	1	8	1	1	4.0	6.0	0	1	2	1
479	210000000.0	4	139.0	0	4	1	1	4.0	6.0	0	1	0	1
399	55000000.0	3	91.0	1	5	1	0	4.0	6.0	0	1	2	1
434	130000000.0	3	130.0	7	5	1	1	3.0	12.0	1	1	2	1

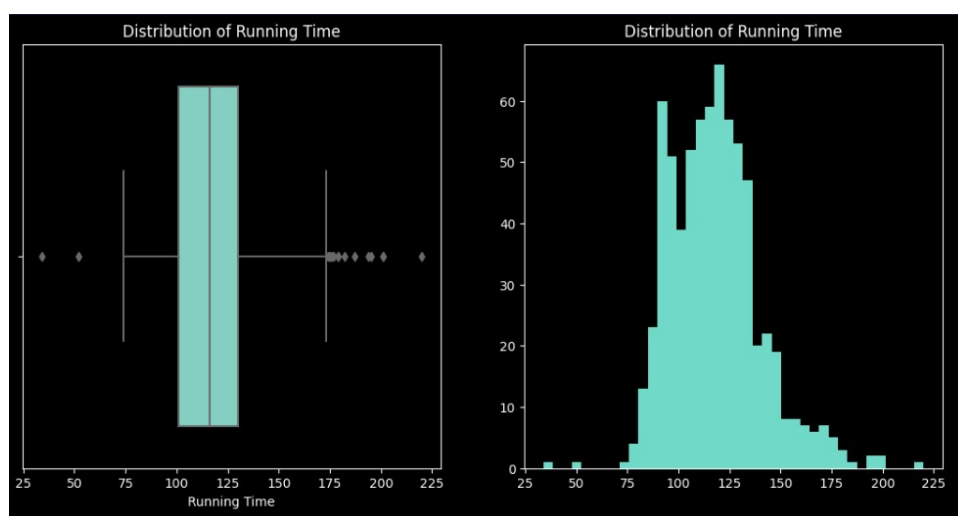
Hình 12. Dữ liệu huấn luyện sau khi được mã hóa

### 3.4 Xử lý ngoại lệ

- Sau khi chia bộ dữ liệu thành 2 tập huấn luyện và kiểm thử với kích thước tập kiểm thử là 30%. Ta vẽ biểu đồ boxplot để kiểm tra ngoại lệ của 2 đặc trưng trong tập dữ liệu huấn luyện.



Hình 13. Phân bố dữ liệu huấn luyện của đặc trưng *Production Budget* khi chưa xử lý ngoại lệ



Hình 14. Phân bố dữ liệu huấn luyện của đặc trưng *Running Time* khi chưa xử lý ngoại lệ

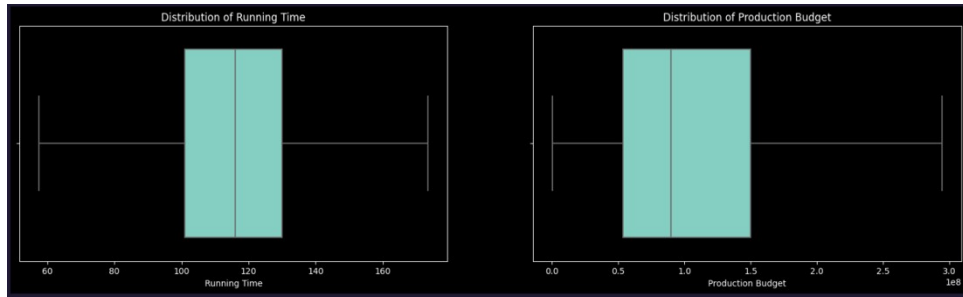
- Hai đặc trưng trên đều có phân bố lệch, nên ta sẽ sử dụng IQR để tìm biên cho phần xử lý ngoại lệ.

- Xử lý ngoại lệ với phân bố lệch:

$$+ \text{Biên trên} = \text{quantile}(0.75) + 1.5 * \text{iqr}$$

$$+ \text{Biên dưới} = \text{quantile}(0.25) - 1.5 * \text{iqr}$$

- Kết quả phân bố dữ liệu sau khi xử lý ngoại lệ được thể hiện qua biểu đồ sau.



Hình 15. Phân bố dữ liệu huấn luyện sau khi xử lý ngoại lệ

### 3.5 Chuẩn hóa dữ liệu

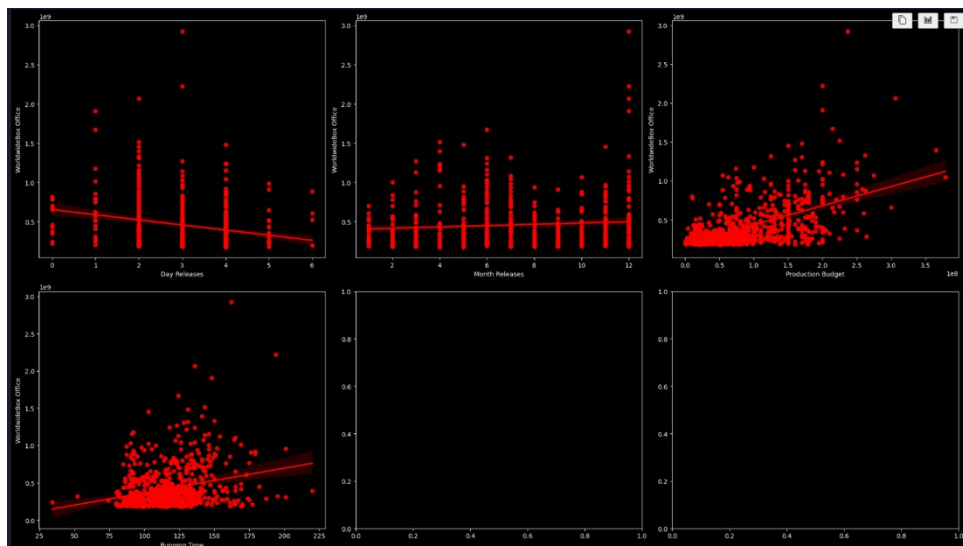
- Chuẩn hóa dữ liệu sử dụng StandardScaler của thư viện sklearn với thuộc tính fit\_transform cho tập dữ liệu huấn luyện và tập dữ liệu kiểm thử.

	Production Budget	MPAA	Running Time	Genre	Creative Type	topCompanies	topDirector	Day Releases	Month Releases	holiday_month	origin_isEnglish	topFranchise	countries_isUS
105	580000000.0	7	107.0	0	8	1	1	1.0	2.0	0	1	2	1
68	258000000.0	4	139.0	1	8	1	1	4.0	6.0	1	1	2	1
479	210000000.0	4	139.0	0	4	1	1	4.0	6.0	1	1	0	1
399	350000000.0	3	91.0	1	5	1	0	4.0	6.0	1	1	2	1
434	130000000.0	3	130.0	7	5	1	1	3.0	12.0	1	1	2	1

Hình 16. Dữ liệu huấn luyện sau khi chuẩn hóa dữ liệu

### 3.6 Độ tương quan giữa các đặc trưng với biến mục tiêu

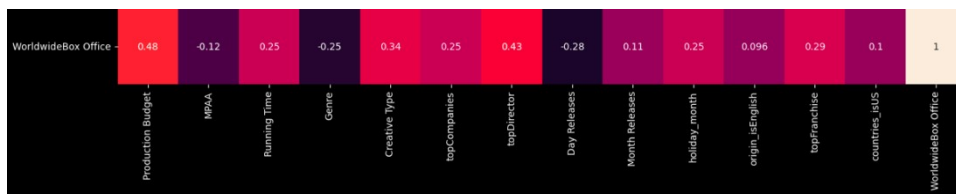
- Tiến hành trực quan hóa sự tương quan của các đặc trưng dạng số so với đặc trưng mục tiêu là 'WorldwideBox Office' bằng sơ đồ regplot. Ta thu được sơ đồ dưới đây:



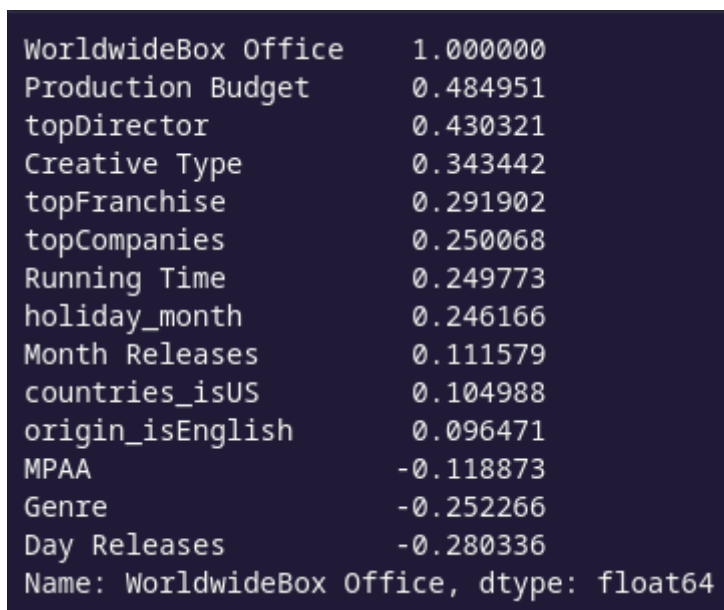
Hình 17. Sự tương quan giữa các đặc trưng với biến mục tiêu

- Có thể đặc trưng 'Production Budget', 'Running Time' có độ tương quan khá cao. Khi giá trị của các đặc trưng này tăng thì doanh thu phim cũng tăng. Ngoài ra, đặc trưng 'Day

Releases' cũng có sự tương quan nhẹ. Còn đặc trưng 'Month Releases' có độ tương quan không cao.



Hình 18. Sự tương quan giữa các đặc trưng với biến mục tiêu



Hình 19. Độ tương quan giảm dần của các đặc trưng so với biến mục tiêu

- Có thể thấy 'Production Budget' có độ tương quan với doanh thu cao nhất. Bộ dữ liệu có các đặc trưng với mức độ tương quan khá thấp cao so với đặc trưng mục tiêu là 'WorldwideBox Office'- Doanh thu phim
- Tổng cộng có 12 đặc trưng được lựa chọn cho mô hình dự đoán: 'Production Budget', 'MPAA', 'Running Time', 'Genre', 'Creative Type', 'topCompanies', 'Top Director', 'Day Release', 'Month Release', 'holiday\_month', 'origin\_isEnglish', 'topFranchise', 'countries\_isUS'.

## 4. Mô hình hóa dữ liệu

### 4.1 Các metrics dùng để đánh giá:

a) Độ chính xác (accuracy):

- Tỷ lệ các trường hợp dự đoán đúng trên tổng số các trường hợp.

$$Accuracy = \frac{TP + TN}{total\ sample}$$

b) Precision:

- Tỷ lệ dự đoán positive đúng trên tổng số trường hợp positive dự đoán.

$$Precision = \frac{TP}{total\ predicted\ positive} = \frac{TP}{TP + FP}$$

c) Recall:

- Tỷ lệ dự đoán positive đúng trên tổng số trường hợp positive thực.

$$Recall = \frac{TP}{total\ actual\ positive} = \frac{TP}{TP + FN}$$

d) F1 Score:

- Trung bình điều hòa giữa precision và recall. Do đó nó đại diện hơn trong việc đánh giá độ chính xác trên đồng thời precision và recall.

$$F_1 = \frac{2}{precision^{-1} + recall^{-1}}$$

e) Các chỉ số TP, FP, TN, FN lần lượt có ý nghĩa là:

- TP (True Positive): Tổng số trường hợp dự báo khớp Positive.
- TN (True Negative): Tổng số trường hợp dự báo khớp Negative.
- FP (False Positive): Tổng số trường hợp dự báo các quan sát thuộc nhãn Negative thành Positive.
- FN (False Negative): Tổng số trường hợp dự báo các quan sát thuộc nhãn Positive thành Negative.

Positive và Negative được qui ước là 1 và 0

## 4.2 Mô hình Logistic Regression:

### a) Cơ sở lý thuyết:

Hồi quy logistic là một kỹ thuật data analysis thường là toán học để xác định sự liên quan giữa hai yếu tố của data. Sau đó, dựa vào sự liên quan này để tiến hành đưa ra các dự đoán của 1 yếu tố khi cho biết yếu tố còn lại.

Giống như tất cả các phân tích Regression, Regression logistic là một phân tích dự đoán. Regression logistic được sử dụng để mô tả dữ liệu và giải thích mối quan hệ giữa một biến nhị phân phụ thuộc và một hoặc nhiều biến độc lập cấp danh nghĩa, thứ tự, khoảng hoặc tỷ lệ.

Hồi quy logistic là một kỹ thuật quan trọng trong lĩnh vực trí tuệ nhân tạo và máy học (AI/ML). Mô hình ML là các chương trình phần mềm có thể được đào tạo để thực hiện các tác vụ xử lý dữ liệu phức tạp mà không cần sự can thiệp của con người. Mô hình ML được xây dựng bằng hồi quy logistic có thể giúp các tổ chức thu được thông tin chuyên sâu hữu ích từ dữ liệu kinh doanh của mình. Họ có thể sử dụng những thông tin chuyên sâu này để phân tích dự đoán nhằm giảm chi phí hoạt động, tăng độ hiệu quả và đổi mới quy mô nhanh hơn.

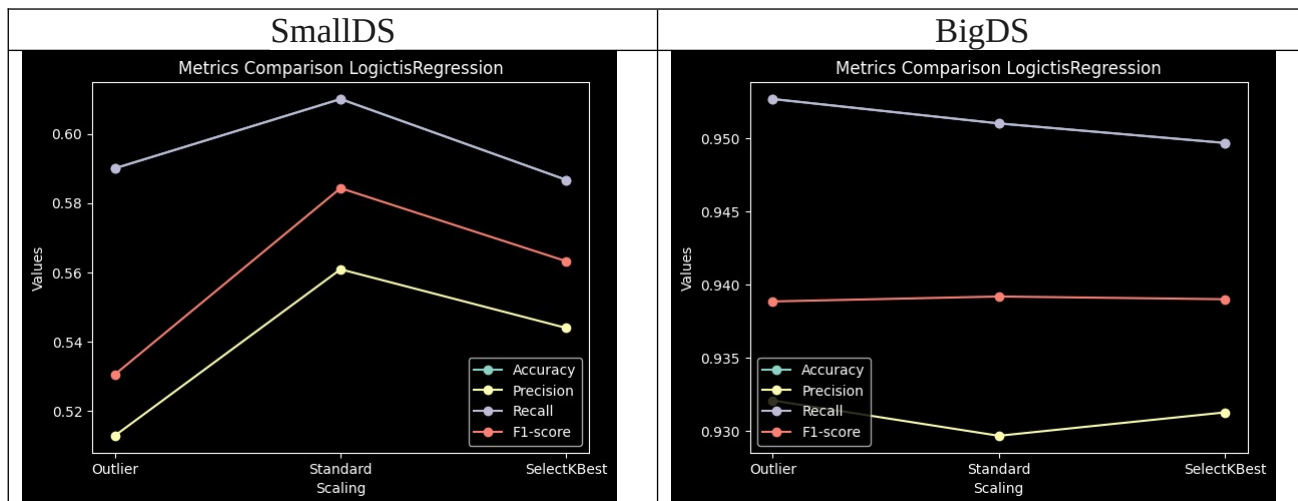
### b) Chia dữ liệu và train mô hình

- Bộ tham số của mô hình: Sử dụng thuật toán RandomizedSearchCV để tìm bộ siêu tham số tối ưu cho mô hình.

Tên tham số	Mảng giá trị	Ý nghĩa
C	uniform(loc=0, scale=4)	Tham số này quy định độ mạnh mẽ của regularization. Càng nhỏ, mô hình càng được regularization mạnh hơn và ngược lại.
penalty	['l1', 'l2']	Tham số này quy định phương pháp sử dụng để xử lý đa cộng tuyến
solver	['liblinear', 'saga']	Tham số này xác định phương pháp sử dụng để tìm nghiệm của mô hình

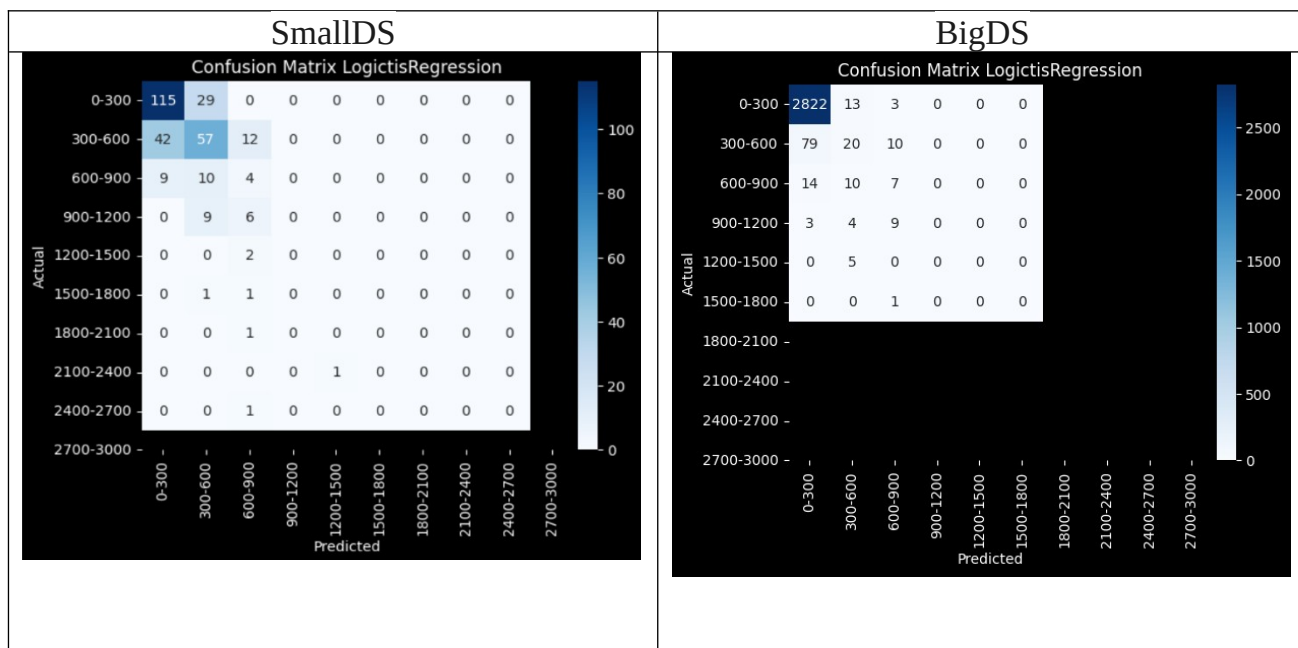
- Chia tập dữ liệu Huấn luyện/Kiểm thử: theo mô hình 30% kiểm thử.

c) Đồ thị thể hiện hiệu suất



Hình 20. Đồ thị thể hiện hiệu suất 4 metrics trên mô hình logistic regression

d) Đồ thị dự đoán so với thực tế



Hình 21. Ma trận nhầm lẫn kết quả dự đoán với kết quả thực của test của mô hình Logistic Regression

e) Đánh giá mô hình theo các metrics là Accuracy, Precision, Recall, F1-Score

SmallDS					BigDS				
	Accuracy	Precision	Recall	F1-score		Accuracy	Precision	Recall	F1-score
Outlier	0.593333	0.513710	0.593333	0.532991	Outlier	0.952667	0.932077	0.952667	0.938840
Standard	0.606667	0.581959	0.606667	0.592468	Standard	0.951000	0.929664	0.951000	0.939178
SelectKBest	0.613333	0.551845	0.613333	0.577135	SelectKBest	0.949667	0.931263	0.949667	0.938991

Hình 22. Đánh giá mô hình Logistic Regression theo các metrics là Accuracy, Precision, Recall, F1-Score

### 4.3 Mô hình Random Forest

#### a) Cơ sở lý thuyết

Random forest là thuật toán supervised learning, có thể giải quyết cả bài toán regression và classification

Random là ngẫu nhiên, Forest là rừng, nên ở thuật toán Random Forest mình sẽ xây dựng nhiều cây quyết định bằng thuật toán Decision Tree, tuy nhiên mỗi cây quyết định sẽ khác nhau (có yếu tố random). Sau đó kết quả dự đoán được tổng hợp từ các cây quyết định.

#### b) Chia dữ liệu và train mô hình

- Bộ tham số của mô hình: Sử dụng thuật toán RandomizedSearchCV để tìm bộ siêu tham số tối ưu cho mô hình.

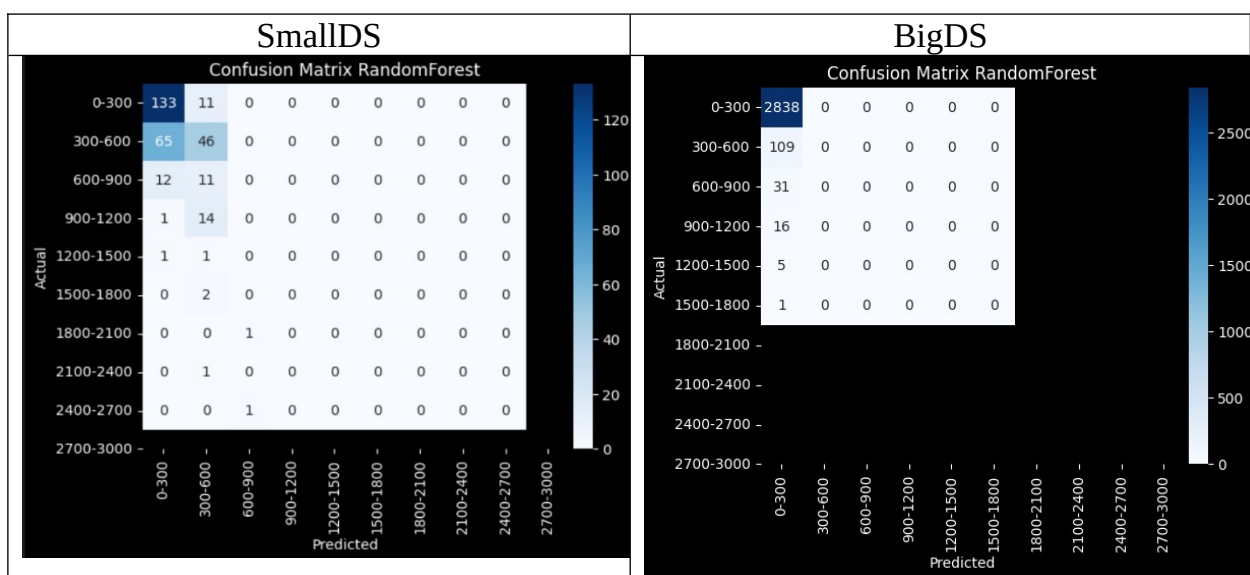
Tên tham số	Mảng giá trị	Ý nghĩa
n_estimators	[10, 50, 100]	Số lượng cây trong rừng (forest)
max_depth	[None, 5, 10]	Độ sâu tối đa của các cây trong rừng. Nếu không được thiết lập, các cây sẽ được phát triển cho đến khi tất cả các lá đều thuần nhất hoặc số mẫu tối thiểu cần thiết để chia tiếp tục không thể đạt được.
min_samples_split	uniform(loc=0, scale=1)	Số lượng mẫu tối thiểu cần thiết để chia một nút trong cây. Nếu một nút



		có số lượng mẫu ít hơn giá trị này, nó sẽ không được chia
min_samples_leaf	uniform(loc=0, scale=0.5)	Số lượng mẫu tối thiểu cần thiết để tạo một lá trong cây. Nếu một lá có số lượng mẫu ít hơn giá trị này, nó sẽ không được tạo

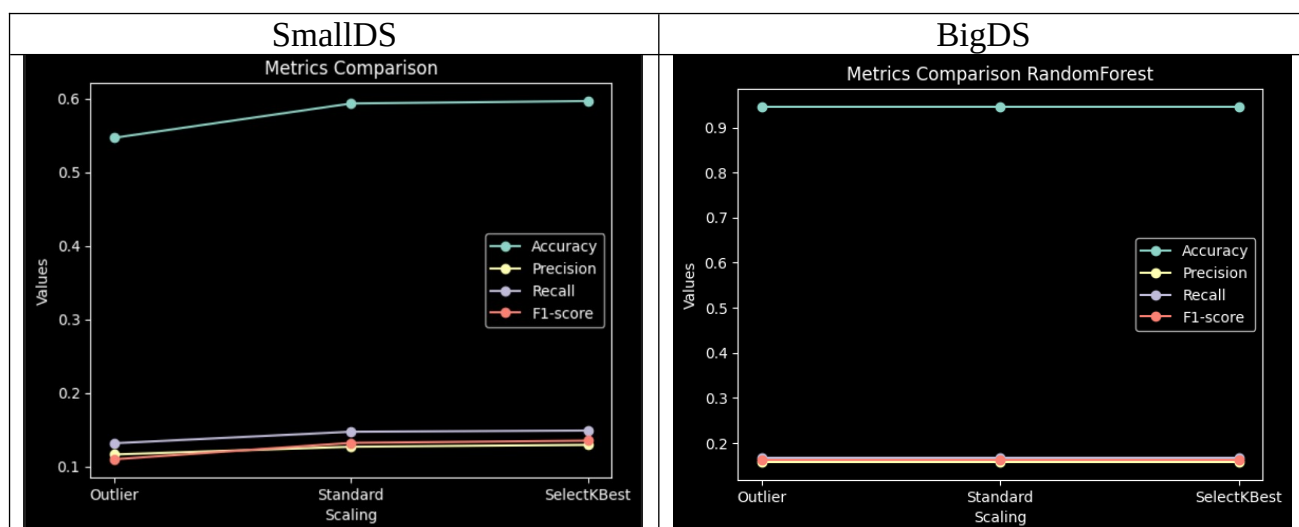
- Chia tập dữ liệu Huấn luyện/Kiểm thử: theo mô hình 30% kiểm thử.

c) Độ lệch dự đoán so với thực tế



Hình 23. Ma trận nhầm lẫn kết quả dự đoán với kết quả thực của test của mô hình Random Forest

d) Đồ thị thể hiện hiệu suất



Hình 24. Đồ thị thể hiện hiệu suất của 4 metrics trên mô hình Random Forest

e) Đánh giá mô hình theo các metrics là Accuracy, Precision, Recall, F1-Score

SmallDS					BigDS				
	Accuracy	Precision	Recall	F1-score		Accuracy	Precision	Recall	F1-score
Outlier	0.570000	0.119874	0.139285	0.121892	Outlier	0.946	0.157667	0.166667	0.162042
Standard	0.576667	0.123839	0.140828	0.123154	Standard	0.946	0.157667	0.166667	0.162042
SelectKBest	0.570000	0.119444	0.138597	0.119863	SelectKBest	0.946	0.157667	0.166667	0.162042

Hình 25. Đánh giá mô hình theo các metrics là Accuracy, Precision, Recall, F1-Score theo Random Forest

## 5. Kết luận

### Những việc đã làm và kết quả:

- Đã tải được dữ liệu từ web.
- Trích xuất và chuẩn hoá dữ liệu.
- Xây dựng và kiểm thử hai mô hình.
- Đánh giá doanh thu dự đoán được từ hai mô hình Logistic Regression và Random Forest:
  - Với SmallDS cả hai mô hình cho kết quả dự đoán trung bình.
  - Với BigDS cả hai mô hình cho ra kết quả dự đoán chính xác hơn rất nhiều.
  - Mô hình Logistic Regression sẽ phù hợp hơn so với Random Forest.

### Hướng phát triển:

- Nếu thêm đặt trưng Reviews, Score, Rating, Population thì sẽ tăng mạnh độ chính xác của mô hình, vì các bộ phim điểm cao thì sẽ có chất lượng tốt và được xem nhiều hơn. Tuy nhiên nếu áp dụng đặt trưng này vào mô hình thì sau này x\_predict cũng sẽ cần đặt trưng các đặc trưng Reviews, Score, Rating, Population, đồng nghĩa với việc bộ phim đã phải hoàn thành đánh giá điểm số (đi ngược lại với mục đích ban đầu của nhóm, đó là dự đoán khi bộ phim còn trong giai đoạn lên ý tưởng).

## 6. Tài liệu tham khảo

- [1] [Logistic Regression - Bài toán cơ bản trong Machine Learning \(viblo.asia\)](#)
- [2] [What is Logistic regression? | IBM](#)
- [3] [Random Forest algorithm — Machine Learning cho dữ liệu dạng bảng \(machinelearningcoban.com\)](#)
- [4] [What is Random Forest? | IBM](#)
- [5] [Web Crawling in Python - MachineLearningMastery.com](#)