

Phân tích phim để dự đoán khả năng thương mại của chúng đối với nhà sản xuất

Devendra Swami

Đại học Nam California Los Angeles, California
dswami@usc.edu

Aadiraj Batlaw

Đại học California Berkeley,
California batlaw33375@berkeley.edu

Yash Phogat

Đại học Nam California Los Angeles, California
phogat@usc.edu

Ashwin Goyal

Đại học Nam California Los Angeles, California
ashwingo@usc.edu

TÓM TẮT

Khi bộ phim ra mắt, một hình thức suy đoán chính liên quan đến thành công thương mại của bộ phim. Tính thương mại này đặc biệt liên quan đến ngân sách ban đầu của bộ phim - vì nhiều lần "phim bom tấn kinh phí lớn" đã đạt được thành công đặc biệt cũng như thất bại thảm hại. Vậy làm thế nào để dự đoán được thành công của một bộ phim sắp ra mắt? Trong bài báo này, chúng tôi đã khám phá một mảng dữ liệu phim lớn trong nỗ lực phát triển một mô hình có thể dự đoán lợi nhuận dự kiến của một bộ phim sắp ra mắt. Cách tiếp cận để phát triển này như sau: Đầu tiên, chúng tôi bắt đầu với tập dữ liệu MovieLens [2] có các thuộc tính phim chung cùng với thể loại cho mỗi phim. Thể loại cung cấp thông tin chi tiết về những đặc điểm cụ thể nào của bộ phim là nổi bật nhất. Sau đó, chúng tôi đã đưa vào các tính năng bổ sung liên quan đến nội dung phim, dàn diễn viên/đoàn làm phim, nhận thức của khán giả, ngân sách và thu nhập từ các trang web TMDB, IMDb và Metacritic. Tiếp theo, chúng tôi đã thực hiện phân tích dữ liệu khám phá và thiết kế một loạt các tính năng mới để nắm bắt thông tin lịch sử cho các tính năng có sẵn. Sau đó, chúng tôi đã sử dụng phân tích giá trị kỳ vọng (SVD) để giảm chiều của các tính năng có chiều cao (ví dụ: thể loại gen). Cuối cùng, chúng tôi đã xây dựng một Bộ phân loại Rừng ngẫu nhiên và thực hiện điều chỉnh siêu tham số để tối ưu hóa độ chính xác của mô hình. Ứng dụng trong thương mại của mô hình của chúng tôi có thể được thấy trong ngành công nghiệp phim ảnh, cho phép các công ty sản xuất dự đoán tốt hơn lợi nhuận dự kiến của các dự án dựa trên phác thảo về quy trình sản xuất của họ, do đó cho phép họ sửa đổi kế hoạch để đạt

- Chúng tôi đã nghiên cứu nhiều đặc điểm có khả năng liên quan đến thành công thương mại của một bộ phim. Khác với các nghiên cứu có sẵn khác, chúng tôi đã kết hợp nhiều đặc điểm mới lạ như công khai, ngày phát hành và dàn diễn viên & đoàn làm phim. Chúng tôi đã dành nhiều thời gian vào kỹ thuật tính năng để hiểu rõ hơn những yếu tố nào khiến một bộ phim có lợi nhuận về mặt tài chính.
- Chúng tôi đã trích xuất 11 nhóm tính năng khác nhau và xây dựng một mô hình rừng ngẫu nhiên (RF) để dự đoán lợi tức đầu tư (ROI) cho bộ phim sẽ cao hơn hay thấp hơn mức trung bình. Sau khi đào tạo RF, chúng tôi đã xác định được tầm quan trọng thương mại của từng tính năng riêng lẻ và các nhóm tính năng.

Bài báo còn lại được cấu trúc như sau. Trong Phần 2, chúng tôi mô tả phương pháp nghiên cứu của mình. Trong Phần 3, chúng tôi trình bày những phát hiện của mình, tiếp theo là Phần 4, nơi chúng tôi thảo luận về mối quan hệ giữa một số tính năng quan trọng và ROI. Trong Phần 5, chúng tôi xác định các mối đe dọa đối với tính hợp lệ và cuối cùng trong Phần 6, chúng tôi kết luận bài báo và đề xuất phạm vi cho công việc trong thương mại.

2 PHƯƠNG PHÁP NGHIÊN CỨU

2.1 Phát biểu vấn đề

Bài báo hiện tại nghiên cứu ứng dụng của việc sử dụng thuật toán học máy để cung cấp hiểu biết tốt hơn về các tính năng có khả năng ảnh hưởng đến thành công thương mại của một bộ phim. Phát biểu vấn đề có thể được định nghĩa chính thức như sau Nhiệm vụ: Dự đoán thành công của bộ phim: Cho một bộ phim

, dự đoán

liệu đây có phải là một bộ phim thành công về mặt thương mại hay không.

Trong khi thực hiện nhiệm vụ này, chúng tôi cố gắng trả lời các câu hỏi nghiên cứu để cập nhật đây.

- RQ1: Thuật toán rừng ngẫu nhiên có thành công như thế nào trong việc dự đoán liệu một bộ phim có thành công về mặt thương mại xét về mặt ROI hay không?
- RQ2: Những tính năng và nhóm tính năng riêng lẻ nào đóng vai trò quan trọng nhất trong việc dự đoán ROI từ phim ảnh?

2.2 Tính năng

Trong phần này, chúng tôi thảo luận về các tính năng mà chúng tôi đã xem xét trong nghiên cứu này. Các tính năng được lựa chọn dựa trên các tính năng thương mại được sử dụng để thực hiện phân tích ngành công nghiệp phim ảnh và các tính năng khác

202411062101.01697v1

1 GIỚI THIỆU

Ngành công nghiệp điện ảnh là một ngành nổi bật ở mọi khía cạnh. Đây là một thế giới riêng. Đối với bài báo này, động lực của chúng tôi xuất phát từ mong muốn cung cấp một mô hình dự đoán cho các nhà sản xuất để có được ý tưởng về khả năng thương mại của bộ phim mà họ đề xuất. Joe Swanberg đã nói vào năm 2016, "Cách duy nhất để bạn kiểm được tiền là nếu bạn đầu tư vào các bộ phim của chính mình". Trước khi các nhà sản xuất phim phải hoàn tất quyết định của mình, họ phải đảm bảo rằng khoản đầu tư của họ là hợp lý và hiểu cách họ thấy được lợi nhuận từ khoản đầu tư đó - đây là nơi mô hình của chúng tôi bước vào thế giới điện ảnh.

Những đóng góp chính trong việc giải quyết vấn đề này có thể được tóm tắt như sau:

Giải nhì, Khu vực Bờ Tây, Citadel Securities, tháng 10 năm 2020, CA, Hoa Kỳ

các tính năng mới có sẵn công khai và có thể dễ dàng trích xuất bằng các công cụ có thể truy cập được.

Bảng 1 tóm tắt các tính năng được xem xét trong nghiên cứu này.

Chúng tôi phân loại các tính năng đã chọn thành 11 nhóm dựa trên đặc điểm của một tính năng trong nhóm. Mỗi nhóm tính năng được tóm tắt dưới đây.

2.2.1 Nội dung. Trong nhóm này, chúng tôi xem xét `is_adult`, `is_english`, `languages_count`, `runtime`, `genome` và `genre`. Nội dung của phim rất quan trọng vì người xem phim thường có sở thích riêng biệt đối với các danh mục này (ví dụ: người nói tiếng Anh thường thích phim tiếng Anh), do đó nội dung phim cho phép chúng tôi thu hẹp đối tượng mục tiêu.

2.2.2 Quảng cáo. Một số tính năng đã được trích xuất để đo lường các nỗ lực quảng cáo, như `is_collection`, `is_homepage`, `is_tagline` và `key-words_count`. Quảng cáo là một yếu tố quan trọng, vì tiếp thị và khả năng tiếp cận của bộ phim có tác động trực tiếp đến số lượng người nghe về bộ phim và theo đó là số lượng người xem bộ phim (đặc biệt là ngày nay, các bộ phim thuộc loại phim có xu hướng vượt trội hơn các bộ phim khác về doanh thu phòng vé).

2.2.3 Nhận thức của khán giả. Trong nhóm này, chúng tôi đã xem xét các tính năng bao gồm mức độ phổ biến, `vote_average`, `vote_count`, `metacritic_score`, `imdb_rating`, `imdb_votes`. Đương nhiên, nhận thức của khán giả có tác động rất lớn đến thành công của phim. Những bộ phim có xếp hạng cao có xu hướng vượt trội hơn những bộ phim có xếp hạng thấp hơn. Tuy nhiên, vì thông tin này thường không có sẵn cho đến sau khi bộ phim được phát hành, nên chúng tôi đã không đưa thông tin này vào mô hình dự đoán của mình.

2.2.4 Ngày phát hành. Trong nhóm này, chúng tôi đã kiểm tra các tính năng như `release_month`, `movies_per_month`, `budget_fraction`, `expense_score`. Một số nghiên cứu đã chỉ ra rằng sự cạnh tranh mà một bộ phim gặp phải tại thời điểm phát hành có tác động lớn đến thành công về mặt tài chính của bộ phim. Nhóm tính năng này cố gắng xem xét hiệu ứng này.

2.2.5 Tài chính. Nhóm này chỉ bao gồm các tính năng, `time_discounted_budget`, là giá trị chiết khấu của ngân sách - được sử dụng để sản xuất phim. Thông thường, những bộ phim có kinh phí lớn hơn có xu hướng vượt trội hơn những bộ phim có kinh phí thấp hơn về mặt doanh thu.

2.2.6 Nhà sản xuất. Nhóm này chỉ bao gồm tính năng `production_house_embedding`. Tính năng này được tính toán từ hiệu suất trung bình của các bộ phim gần đây do cùng một nhà sản xuất sản xuất. Các hãng sản xuất phim tốt và uy tín hơn thường có thể thuê được nhiều đạo diễn, biên kịch và ngôi sao nổi tiếng hơn, ngoài việc có ngân sách lớn hơn, đây cũng là những yếu tố thường dẫn đến thành công về mặt doanh thu phòng vé.

2.2.7 Biên kịch. Nhóm này chỉ bao gồm tính năng `writers_embedding`, được tính toán từ hiệu suất trung bình của các bộ phim gần đây do cùng một biên kịch viết. Thông thường, các biên kịch đã viết cho những bộ phim thành công có xu hướng viết thành công hơn

phim ảnh.

2.2.8 Đạo diễn. Nhóm này chỉ bao gồm các đạo diễn đặc trưng `embedding`, được tính toán từ hiệu suất trung bình của các bộ phim gần đây do cùng một đạo diễn thực hiện. Thông thường, các đạo diễn đã chỉ đạo các bộ phim thành công có xu hướng chỉ đạo các bộ phim thành công hơn

phim ảnh.

2.2.9 Nhà sản xuất. Nhóm này chỉ bao gồm các tính năng `producers_embedding`, được tính toán từ hiệu suất trung bình của các bộ phim gần đây do cùng một nhà sản xuất thực hiện. Thông thường, các nhà sản xuất đã sản xuất những bộ phim thành công có xu hướng sản xuất những bộ phim thành công hơn.

2.2.10 Diễn viên chính. Nhóm này chỉ bao gồm tính năng `main_cast_embedding`, được tính toán từ hiệu suất trung bình của các bộ phim gần đây mà các thành viên từ dàn diễn viên chính tham gia

Nhiều nghiên cứu đã chỉ ra rằng ngôi sao của bộ phim có ảnh hưởng lớn đến sự thành công của bộ phim.

2.2.11 Nhân viên hỗ trợ. Nhóm này bao gồm các tính năng `fe-male_count`, `male_count` và `crew_length`. Một số nghiên cứu chỉ ra rằng giới tính và quy mô của dàn diễn viên trong một bộ phim có thể ảnh hưởng đến thành công.

2.3 Thu thập dữ liệu

Chúng tôi đã lấy được hầu hết dữ liệu từ tập dữ liệu Kaggle "The Movies Dataset".[3] cung cấp cho chúng tôi dữ liệu về dàn diễn viên, đoàn làm phim và siêu dữ liệu chung trên một tập hợp con lớn các bộ phim. Trong một tập dữ liệu do cộng đồng cuộc thi Data Open cung cấp, chúng tôi đã lấy được các thẻ bộ gen, chúng tôi đã hợp nhất chúng với tập dữ liệu siêu dữ liệu của mình. Chúng tôi đã lấy thêm các tính năng bằng cách thực hiện các lệnh gọi API đến TMDB (The Movie Database), cũng như thu thập điểm Metacritic từ "IMDB: All US Phim đã phát hành: 1972-2016". Ban đầu, tập dữ liệu hợp nhất của chúng tôi bao gồm hơn 13 nghìn hàng một chút vì chúng tôi chỉ có thông tin bộ gen cho nhiều bộ phim này. Sau đó, chúng tôi quyết định tập trung vào các bộ phim sau cuộc cách mạng màu (tức là sau năm 1965) và loại bỏ các hàng có giá trị ngân sách hoặc giá trị doanh thu được đặt thành 0. (Tuy nhiên, chúng tôi đã cân nhắc trước để nội suy các lỗi này, vì ngân sách và doanh thu là một yếu tố quan trọng trong việc xác định tỷ lệ hoàn vốn, nên chúng tôi kết luận rằng bất kỳ phương pháp nội suy nào cũng sẽ làm sai lệch mô hình của chúng tôi phần lớn). Điều này đưa tổng số hàng của chúng tôi lên 5.426.

2.4 Thuật toán học máy Lúc đầu chúng tôi nghĩ đến việc

sử dụng hồi quy làm thuật toán học máy của mình do thực tế là các giá trị ROI là liên tục. Tuy nhiên, trong nhiều trường hợp, dự đoán của chúng tôi từ hồi quy không chính xác. Hơn nữa, thay vì quan tâm đến việc biết các giá trị chính xác, các nhà sản xuất phim có thể có xu hướng muốn biết liệu bộ phim của họ có khả năng hoạt động tốt hay không. Do đó, chúng tôi đã quyết định dự đoán trên hoặc dưới ROI trung bình (trung bình được tính trên dữ liệu đào tạo) thay vì các giá trị ROI chính xác. Theo cách này, chúng tôi đã chuyển đổi vấn đề hồi quy khó học ban đầu của mình thành một nhiệm vụ phân loại nhị phân dễ dàng.

Chúng tôi đã triển khai thuật toán rừng ngẫu nhiên (RF) để thực hiện nhiệm vụ phân loại của mình vì đây là một trong những thuật toán phân loại phi tuyến tính thành công nhất.

Thuật toán học máy. Nó xem xét nhiều cây quyết định được đào tạo trên mẫu ngẫu nhiên của dữ liệu đào tạo. Ngoài ra, để phân chia các nút, nó chọn ngẫu nhiên một tập hợp con các tính năng. Cuối cùng, quyết định được thực hiện bằng cách lấy trung bình các dự đoán từ mỗi cây quyết định. Do đó, thuật toán rừng ngẫu nhiên ít có khả năng bị quá khớp. Hơn nữa, nó có thể lấy các biến số hoặc biến danh mục làm đầu vào và thậm chí không yêu cầu mở rộng tính năng. Vì những lý do này, chúng tôi đã chọn thuật toán rừng ngẫu nhiên để thực hiện nhiệm vụ phân loại. Tập dữ liệu được chia thành khoảng 80% đào tạo

Phân tích phim để dự đoán khả năng thương mại của chúng đối với nhà sản xuất

Giải nhì, Khu vực Bờ Tây, Citadel Securities, tháng 10 năm 2020, CA, Hoa Kỳ

Bảng 1: Các tính năng của phim có khả năng ảnh hưởng đến thành công thương mại của một bộ phim

| # | Nhóm | Tính năng | Sự miêu tả |
|----|------------------------|---------------------------------|--|
| 1 | Nội dung | là_người_đi_lớn | Bộ phim này có chỉ phù hợp với người đi lớn hay không. |
| | | là_tiếng_anh | Ngôn ngữ chính được sử dụng trong phim có phải là tiếng Anh hay không. |
| | | ngôn | Tổng số ngôn ngữ được sử dụng trong phim. |
| | | ngữ_số_lượng_phim_runtime | Tổng độ dài của bộ phim, tính bằng phút. |
| | | thể_genome [6] | Các tính năng được trích xuất từ nội dung phim. |
| 2 | Công khai | thể_loại_phim | Những thể loại khác nhau mà bộ phim có thể được xếp vào. |
| | | là_bộ_sưu_cung | Bộ phim có phải là một phần của bộ sưu tập hay loạt phim không. |
| | | là_trang_chủ | Bộ phim có trang chủ hay không. |
| | | là_khẩu_hiệu | Bộ phim có khẩu hiệu liên quan hay không. |
| | | số_lượng_từ_khóa | Số lượng từ khóa thứ ởng dùng có thể gắn cho phim. |
| 3 | Nhận thức của khán giả | sự_phổ_biến | Điểm số phổ biến của phim do TMDB cung cấp. |
| | | biên_chọn_trung_bình | Điểm chấp thuận trung bình cho bộ phim, dựa trên tổng số phiếu bầu tại TMDB. |
| | | số_phiếu_bầu | Tổng số phiếu bầu cho bộ phim tại TMDB. |
| | | Điểm_Metacritic | Điểm đánh giá trung bình của phim trên Metacritic. |
| | | Xếp_hạng_imdb | Điểm đánh giá trung bình của phim trên IMDB. |
| 4 | Ngày phát hành | imdb_vote | Tổng số phiếu bầu cho bộ phim trên IMDB. |
| | | tháng_phát_hành | Tháng mà bộ phim được phát hành. |
| | | movies_per_month | Tổng số phim được phát hành trong cùng tháng với bộ phim được chỉ định. |
| | | budget_fraction | Ngân sách của bộ phim tỷ lệ thuận với ngân sách tích lũy của tất cả các bộ phim trong tháng đó. |
| | | điểm_chi_phi_phim | Ngân sách của bộ phim tương ứng với ngân sách trung bình của tất cả các bộ phim trong tháng đó. |
| 5 | Tài chính | time_discounted_budget | Giá trị chiết khấu theo thời gian cho ngân sách dùng để sản xuất phim. |
| 6 | Nhà sản xuất | nhà_sản_xuất_những_writer_những | Tính toán dựa trên hiệu suất của các bộ phim gần đây do cùng một hãng sản xuất sản xuất. |
| 7 | Nhà văn | | Tính toán dựa trên thành tích của những bộ phim gần đây do cùng tác giả chấp bút. |
| 8 | Giám đốc | giám_đốc_những | Tính toán dựa trên thành tích của những bộ phim gần đây do cùng đạo diễn chỉ đạo. |
| 9 | Nhà sản xuất | nhà_sản_xuất_những | Tính toán dựa trên hiệu suất của các bộ phim gần đây do cùng một nhà sản xuất sản xuất. |
| 10 | Diễn viên chính | những_main_cast | Tính toán dựa trên thành tích của các bộ phim gần đây có sự tham gia của các thành viên trong dàn diễn viên chính. |
| 11 | Nhân viên hỗ trợ | nữ_số_lượng | Tổng số diễn viên nữ trong phim. |
| | | nam_số_lượng | Tổng số diễn viên nam trong phim. |
| | | chiều_dài_phil_hành_đoàn | Tổng số người trong đoàn làm phim. |

†: Không được sử dụng trong mô hình dự đoán vì thông tin đó không có sẵn trước khi phim được phát hành.

và 20% dữ liệu thử nghiệm với các bộ phim trước năm 2011 là một phần của chương trình đào tạo dữ liệu và từ năm 2011 trở đi bao gồm bộ thử nghiệm.

Giảm chiều: Do sự hiện diện của các đặc điểm thưa thớt có chiều cao như thể bộ gen và thông tin thể loại (sau một lần mã hóa nóng) trong tập dữ liệu, chúng tôi đã sử dụng giá trị kỳ dị phân tích (SVD) để giảm số chiều của chúng trong khi giữ nguyên hầu hết các phương sai trong dữ liệu. Ngoài ra, chúng tôi cũng đã loại bỏ các tính năng có tương quan cao khỏi tập dữ liệu của chúng tôi. Chúng tôi đã xác định được cặp tính năng có mối tương quan cao, tức là cặp có giá trị tương quan tuyệt đối là 0,75 hoặc cao hơn từ tương quan Spearman . Trong số các cặp được xác định này, các tính năng có mối quan hệ tương hỗ thấp hơn thông tin [4] bị loại bỏ. Nó được thực hiện để loại bỏ các tính năng không cung cấp bất kỳ thông tin bổ sung nào và có thể giúp giảm kích thước của tập dữ liệu để tăng tốc quá trình đào tạo của chúng tôi.

Tối ưu hóa siêu tham số: Để có được siêu tham số tối ưu cho mô hình RF của chúng tôi, chúng tôi đã bắt đầu bằng cách chuẩn bị một không gian tìm kiếm lưới. Đối với n_estimators, chúng tôi đã xem xét các giá trị từ 100 đến 1000 với bước nhảy là 100. Tương tự như vậy, bội số nguyên của √ # đến #features cho max_features, giá trị từ 10 đến 100 với kích thước bước là 10 cho max_depth được sử dụng. Các giá trị 0,01, 0,03,

và 0,05 được xem xét cho min_samples_split, và 1, 3 và 5 cho min_samples_leaf. Chúng tôi luôn đặt bootstrap thành True. Do không gian tìm kiếm lưới của 12600 (10 · 14 · 10 · 3 · 3) là tốn kém về mặt tính toán, chúng tôi đã thực hiện tìm kiếm ngẫu nhiên trên các mẫu được rút ra đồng đều từ nó. Chúng tôi đã sử dụng 100 lần lặp lại thông qua xác thực chéo 4 lần để xác định các siêu tham số tối ưu. Chúng tôi báo cáo các siêu tham số thu được từ quy trình này trong Bảng 2.

Bảng 2: Các tham số của mô hình RF được điều chỉnh trên tập xác thực

| Các tham số | Giá trị tối ưu |
|-------------------|----------------|
| n_estimators | 500 |
| max_features | 14 |
| max_depth | 40 |
| min_samples_split | 0,05 |
| min_samples_leaf | 5 |

2.5 Đánh giá mô hình

Nhờ cách chúng tôi xây dựng nhiệm vụ phân loại của mình, tập dữ liệu của chúng tôi được cân bằng hoàn hảo với hai lớp có tầm quan trọng như nhau. Do đó , Độ chính xác ban đầu được coi là một thước đo đánh giá phù hợp. Tuy nhiên, độ chính xác phụ thuộc vào việc lựa chọn ngưỡng để chuyển đổi xác suất dự đoán thành lớp đầu ra có khả năng xảy ra cao nhất. Khác nhau ngưỡng có thể tạo ra điểm chính xác khác nhau cho cùng một dự đoán ML và do đó gây khó khăn cho việc so sánh hiệu suất trên nhiều thuật toán ML.

Để vượt qua rào cản này, chúng tôi đã chọn Diện tích dưới đường cong Đặc tính hoạt động của Máy thu (ROC) để đánh giá hiệu suất thực hiện nhiệm vụ của chúng tôi. Chúng tôi sử dụng từ viết tắt AUC để biểu thị số liệu này. Đường cong ROC là một biểu đồ giữa tỷ lệ dương thực sự và tỷ lệ dương tính giả thu được bằng cách triển khai các ngưỡng khác nhau để chỉ định các giá trị dương. Diện tích dưới đường cong này (AUC) sau đó được sử dụng để so sánh hiệu suất của một mô hình trên phân loại nhị phân của chúng tôi nhiệm vụ.

Hơn nữa, chúng tôi đã so sánh hiệu suất của thuật toán rừng ngẫu nhiên với phương pháp tiếp cận đường cơ sở ngẫu nhiên. Trong Ngẫu nhiên đường cơ sở, chúng tôi ngẫu nhiên chỉ định dưới hoặc trên ROI trung bình cho mỗi phim trong bộ thử nghiệm. Không có sự cải thiện đáng kể nào trong hiệu suất của mô hình RF trên đường cơ sở sẽ ngụ ý rằng đã chọn các tính năng không liên quan đến thành công thương mại của bộ phim.

Một mô hình phân loại hoàn hảo sẽ mang lại giá trị AUC là 1,00, trong khi một mô hình ngẫu nhiên sẽ có giá trị khoảng 0,50. Nói chung, một giá trị AUC cao hơn cho thấy mô hình tốt hơn.

2.6 Phân tích tầm quan trọng của tính năng

Chúng tôi đã sử dụng kỹ thuật tính năng hoán vị quan trọng được đề xuất bởi Breiman [1] để đo lường tầm quan trọng của một tính năng riêng lẻ và một nhóm tính năng. Ý tưởng chính đằng sau phương pháp này là rằng sự hoán vị ngẫu nhiên của các tính năng quan trọng sẽ dẫn đến một sự suy giảm đáng kể về hiệu suất của mô hình. Sự suy giảm này thước đo hiệu suất được gọi là Giá trị quan trọng (IV) và giá trị càng cao thì tính năng đó càng quan trọng.

Tương tự như vậy, tầm quan trọng của nhóm được tính bằng cách hoán vị tất cả các tính năng trong một nhóm với nhau. Điều quan trọng là không ở đây rằng chúng ta có thể sử dụng an toàn kỹ thuật quan trọng của tính năng hoán vị [5] vì không có tính năng nào có mối tương quan cao trong tập dữ liệu của chúng tôi chúng tôi đã loại bỏ những tính năng có mối tương quan cao.

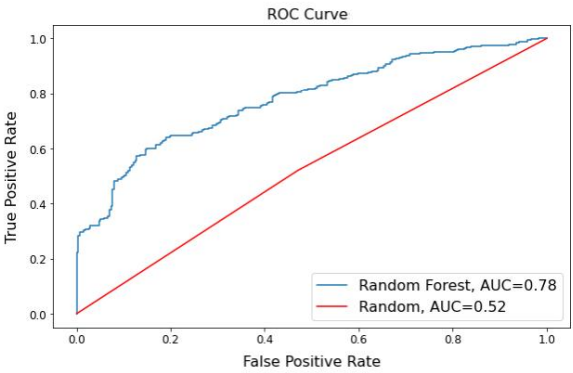
3 KẾT QUẢ

3.1 RQ1: Rừng ngẫu nhiên thành công như thế nào

thuật toán trong việc dự đoán xem một bộ phim sẽ là một thành công thương mại về mặt

Lợi tức đầu tư ?

Hình 1 cho thấy các đường cong ROC của các đường cơ sở và ngẫu nhiên Thuật toán rừng. AUC cho thuật toán rừng ngẫu nhiên là 0,78. Về cơ sở, Random đạt được điểm AUC mong đợi khoảng 0,500. Điểm AUC cao hơn của mô hình RF cho thấy rằng giải pháp sử dụng kỹ thuật tính năng và máy học sẽ có khả năng giúp các nhà sản xuất phim và các nhà phân tích làm tốt hơn ước tính lợi tức đầu tư dự kiến cho các bộ phim sắp ra mắt.



Hình 1: Đường cong ROC của thuật toán rừng ngẫu nhiên sử dụng các tính năng đã chọn và đường cơ sở cho tập dữ liệu của chúng tôi.

3.2 RQ2: Những đặc điểm và nhóm riêng lẻ nào của các tính năng đóng vai trò quan trọng nhất trong dự đoán ROI từ phim ảnh?

Bảng 3: Tầm quan trọng của 15 tính năng cá nhân hàng đầu

| Tên tính năng | Giá trị thông tin (IV) |
|--------------------------------------|------------------------|
| is_collection | 0,010 |
| bộ gen_0 số | 0,005 |
| lưu ứng từ khóa bộ | 0,004 |
| gen_2 bộ | 0,004 |
| gen_13 | 0,004 |
| phim_mỗi_tháng số lưu ứng | 0,003 |
| nam bộ gen_3 | 0,002 |
| bộ gen_5 bộ | 0,002 |
| gen_12 | 0,002 |
| is_homepage | 0,002 |
| thời | 0,001 |
| gian_giảm_giá_ngân_sách_female_count | 0,001 |
| bộ gen_1 bộ gen_10 | 0,001 |
| | 0,001 |
| | 0,001 |

Bảng 3 cho thấy 15 tính năng quan trọng nhất trong việc dự đoán ROI của phim. Chúng tôi cũng quan sát thấy is_collection nằm trong số tính năng quan trọng nhất. Điều này có ý nghĩa trực quan, như phần tiếp theo, các phim phụ và phim mở rộng của một vũ trụ chủ đề có xu hướng vượt trội hơn các phim không có những đặc điểm như vậy. Chúng tôi cũng nhận thấy rằng các đặc điểm của bộ gen (mà như người ta có thể nhớ lại mô tả một bộ phim của (đặc điểm nổi bật nhất của nó) có xu hướng cung cấp thông tin có giá trị. keyword_count cung cấp thông tin tương tự như tính năng bộ gen, và nó cũng có giá trị thông tin cao ngụ ý rằng đặc điểm của phim là một cái nhìn sâu sắc quan trọng để xác định sự thành công của một bộ phim. Quan sát đáng chú ý cuối cùng trong bảng này sẽ là movies_per_month, ngụ ý rằng sự cạnh tranh vào ngày phát hành lại của một bộ phim thực sự là một yếu tố chính góp phần quyết định sự thành công của một bộ phim.

Bảng 4 hiển thị 5 nhóm tính năng quan trọng nhất. Ở đây chúng tôi quan sát rằng Nội dung có giá trị thông tin cao nhất. Điều này làm cho

Phân tích phim để dự đoán khả năng thương mại của chúng đối với nhà sản xuất

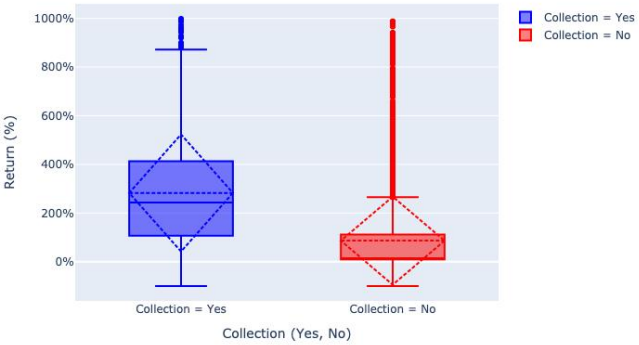
Bảng 4: Tầm quan trọng của 5 nhóm tính năng hàng đầu

| Giá trị thông tin nhóm | tính năng (IV) |
|------------------------|----------------|
| Nội dung | 0,047 |
| Diễn viên chính | 0,031 |
| Công khai | 0,021 |
| Nhà văn | 0,017 |
| Nhà sản xuất | 0,016 |

nghĩa là, vì genome_tags được đưa vào như một tính năng trong danh mục này. Với Main Cast đứng thứ hai trong IV, chúng tôi khẳng định lại giả thuyết của mình rằng các ngôi sao của một bộ phim có tác động lớn đến thành công của bộ phim.

4 THẢO LUẬN

Thay vì chỉ dựa vào các giá trị tầm quan trọng của tính năng đã thảo luận ở phần trước, chúng tôi cũng muốn tìm hiểu mối quan hệ nhân quả và xác định hướng ảnh hưởng của các tính năng quan trọng lên giá trị ROI. Do đó, trong phần này, chúng tôi sẽ khám phá cụ thể mối quan hệ giữa các tính năng chính và lợi nhuận. Tuy nhiên, chúng tôi lưu ý rằng đây là phân tích đơn biến có nhược điểm là không xem xét các tương tác đa biến.



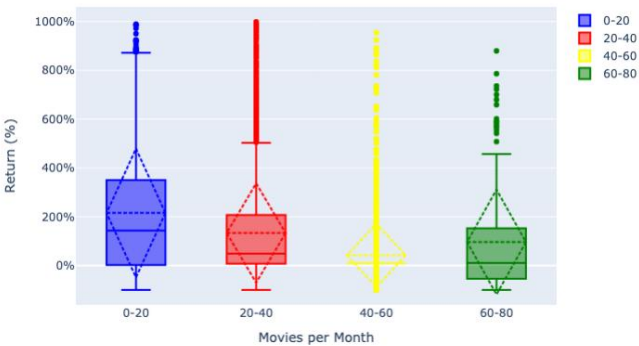
Hình 2: Biểu đồ giữa Is_collection và Return

Hình 2 cho thấy có vẻ như có sự khác biệt hợp lý về lợi tức đầu tư của những bộ phim nằm trong bộ sưu tập và những bộ phim không nằm trong bộ sưu tập. Những bộ phim nằm trong bộ sưu tập/loạt phim có xu hướng mang lại ROI cao hơn những bộ phim không nằm trong bất kỳ bộ sưu tập nào.

Hình 3 cho thấy có vẻ như có sự khác biệt khá hợp lý giữa lợi tức đầu tư cho các bộ phim được phát hành với các phạm vi khác nhau về số lượng phim được phát hành trong tháng đó. Có vẻ như phần lớn, càng ít phim được phát hành trong tháng của ngày phát hành phim quan tâm thì tỷ lệ lợi nhuận càng cao. Điều quan trọng cần lưu ý là có những ngoại lệ lớn cho tất cả các danh mục này, báo hiệu rằng một số phim trên tất cả các danh mục đều có thành tích đặc biệt tốt bất kể số lượng phim được phát hành trong cùng tháng của ngày phát hành phim.

Hình 4 khám phá xem sự tồn tại của trang chủ cho một bộ phim cụ thể có liên quan đến lợi tức đầu tư cao hay không. Từ hình này, có vẻ như có mối liên hệ, nhưng có vẻ như không đặc biệt quan trọng đối với tương ứng hợp đơn biến vì chúng ta thấy các khoảng tin cậy chồng chéo.

Giải nhì, Khu vực Bờ Tây, Citadel Securities, tháng 10 năm 2020, CA, Hoa Kỳ

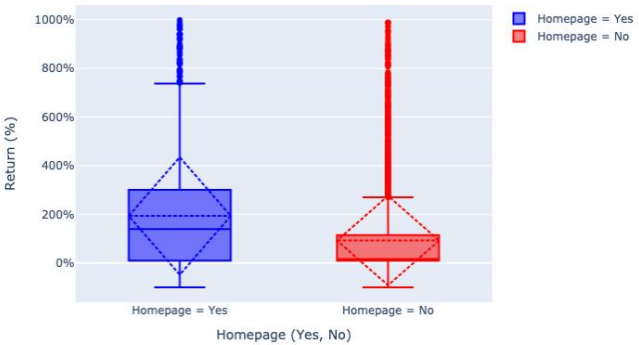


Hình 3: Biểu đồ giữa Phim mỗi tháng và Lợi nhuận

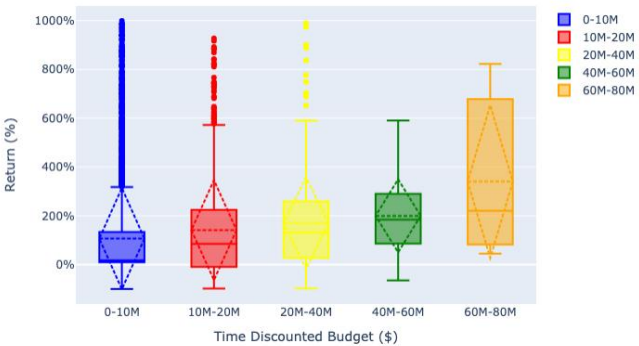
Hình 5 so sánh lợi tức đầu tư cho các loại ngân sách khác nhau. Từ đây chúng ta thấy rằng các bộ phim được sản xuất với ngân sách lớn hơn có xu hướng có lợi tức đầu tư cao hơn.

5 MỐI ĐE DỌA ĐỐI VỚI TÍNH HỢP LỆ

Chúng tôi đã liệt kê nhiều mối đe dọa khác nhau đến tính hợp lệ và hạn chế của nghiên cứu trong phần này.



Hình 4: Biểu đồ giữa Is_homepage và Return



Hình 5: Biểu đồ giữa Ngân sách chiết khấu theo thời gian và Lợi nhuận

Độ lệch lựa chọn tính năng: Kết quả của nghiên cứu hiện tại có thể bị ảnh hưởng bởi các tính năng mà chúng tôi đã xem xét, bị giới hạn bởi trí tư đồng tư đồng của chúng tôi. Do đó, những người khác sử dụng một bộ tính năng khác có thể thu được các kết quả thử nghiệm khác. Có một số tính năng mà chúng tôi đã xem xét để dự đoán nhưng vẫn chưa thể kết hợp, ví dụ nếu bộ phim được đề cử giải thưởng của viện hàn lâm, thì nhóm sản xuất có kết nối tốt như thế nào với thể giới? Có, chúng tôi đã đưa vào các tính năng phổ biến nhưng câu hỏi rộng hơn vẫn là, ngân sách được quản lý như thế nào? Bao nhiêu đã được phân bổ cho tiếp thị và chiến dịch?. Do đó, việc đưa thêm nhiều tính năng hơn sau khi tham khảo ý kiến của các chuyên gia trong ngành là một phần trong kế hoạch tư đồng lai của chúng tôi.

Độ tin cậy của công cụ và phương pháp: Mặc dù chúng tôi đã sử dụng các công cụ và phương pháp chuẩn phù hợp cho các nghiên cứu như vậy, nhưng vẫn có khả năng tồn tại những điểm không chính xác mà chúng tôi không tính đến.

Ví dụ, chúng tôi cho rằng các thể bộ gen được tính toán hoàn hảo và cung cấp biểu diễn tốt hơn về nội dung phim. Tuy nhiên, do được lấy từ thuật toán ML, các thể này có thể không chính xác trong một số trường hợp.

Tính hợp lệ bên ngoài: Chúng tôi đồng ý rằng có một số hạn chế có khả năng ảnh hưởng đến khả năng khái quát hóa các phát hiện của chúng tôi. Đầu tiên, vấn đề lấy mẫu phim; để giải thích rõ hơn, chúng tôi nhận ra rằng thời thể đã thay đổi kể từ đầu những năm 1920 và do đó chúng tôi chỉ cần xem xét một khung thời gian hạn chế trước thời kỳ mục tiêu, điều này dẫn đến việc giảm số điểm dữ liệu và do đó giảm số lần kiểm tra tính hợp lệ mà chúng tôi có thể thực hiện.

Hơn nữa, chỉ có một thuật toán học máy được sử dụng, cụ thể là Random Forest. Chúng tôi đã cân nhắc sử dụng mạng nơ-ron nhưng do thiếu điểm dữ liệu nên mô hình của chúng tôi sẽ bị lỗi.

6 KẾT LUẬN VÀ CÔNG VIỆC TƯ ĐỒNG LAI

Trong bài báo này, chúng tôi tìm cách trả lời câu hỏi các công ty sản xuất có thể dự đoán lợi nhuận trong tư đồng lai của các dự án điện ảnh sắp tới như thế nào, cố gắng giúp các nhà làm phim hiểu rõ hơn về khả năng thương mại của các dự án của họ. Để đạt được điều này, chúng tôi đã trích xuất các tính năng trong các danh mục Nội dung, Quảng cáo, Ngày phát hành, Tài chính, Nhà sản xuất, Biên kịch, Đạo diễn, Nhà sản xuất, Diễn viên chính và Nhân viên hỗ trợ. Cuối cùng, chúng tôi đã xây dựng một mô hình phân loại rừng ngẫu nhiên, sau khi điều chỉnh siêu tham số kỹ lưỡng, có thể dự đoán thành công lợi nhuận của các bộ phim có điểm AUC là

78%.

Công việc trong tư đồng lai: Ưu tiên hàng đầu là thêm nhiều tính năng hơn và mở rộng tập dữ liệu của chúng tôi để bao gồm nhiều điểm dữ liệu hơn cho các tính năng khả dụng. Điều này bao gồm việc thu thập dữ liệu từ các trang web truyền thông xã hội để tìm các kết nối ẩn và mạng lưới sâu trước đây chưa được biết đến, ví dụ, xây dựng biểu đồ kết nối truyền thông xã hội để xác định tốt hơn mức độ phổ biến và các vòng tròn bên trong của nhóm sản xuất, từ đó cung cấp các dự đoán tốt hơn cho các vấn đề cốt lõi của chúng tôi. (Các vòng tròn xã hội xác định mức độ hiện diện của một thành viên duy nhất có sức mạnh như thế nào - sự hiện diện càng tốt, lợi nhuận càng cao). Sau khi thực hiện xong điều đó, chúng tôi muốn tập trung nỗ lực vào việc sử dụng các thuật toán học tập khác để có được các mô hình mới có thể so sánh hoặc kết hợp với mô hình hiện tại của chúng tôi. Một thuật toán như vậy có thể là mạng nơ-ron vì chúng ta có thể khám phá thêm các mối quan hệ trong thế giới thực hiện tại chưa được biết đến. Tất nhiên, điều này có thể thực hiện được khi chúng tôi đã mở rộng tập dữ liệu của mình một cách đầy đủ.

TÀI LIỆU THAM KHẢO

[1] Breiman, L. Rừng ngẫu nhiên. Học máy 45, 1 (2001), 5-32.
[2] Harper, FM và Konstan, JA Bộ dữ liệu movielens: Lịch sử và bối cảnh. ACM Trans. Tư đồng tác. Trí tuệ. Hệ thống. 5, 4 (tháng 12/2015).
[3] Kaggle. Bộ dữ liệu phim. Trực tuyến, 2017.
[4] Kraskov, A., Stögbauer, H., và Grassberger, P. Ước tính thông tin lẫn nhau. Đánh giá vật lý E 69, 6 (2004), 066138.
[5] Strobl, C., Boulesteix, A.-L., Kneib, T., Augustin, T., và Zeileis, A. Tầm quan trọng của biến có điều kiện đối với rừng ngẫu nhiên. BMC bioinformatics 9, 1 (2008), 307.
[6] Vig, J., Sen, S., và Riedl, J. Bộ gen thể: Mã hóa kiến thức cộng đồng để hỗ trợ tư đồng tác mới. ACM Trans. Interact. Intell. Syst. 2, 3 (tháng 9 năm 2012).