
Early Predictions of Movie Success: The Who, What, and When of Profitability

MICHAEL T. LASH AND KANG ZHAO

MICHAEL T. LASH (michael-lash@uiowa.edu) is a Ph.D. student in the Department of Computer Science at University of Iowa. His research interests lie in the areas of data mining, machine learning, and predictive analytics. Specific interests as well as ongoing areas of research include inverse classification, utility-based data mining, adversarial learning, and survival analytics and learning. Application of these areas to health care, business, and entertainment domains are also of interest.

KANG ZHAO (kang-zhao@uiowa.edu; corresponding author) is an assistant professor at Tippie College of Business, University of Iowa. He is also affiliated with the university's Interdisciplinary Graduate Program in Informatics. He obtained his Ph.D. from Penn State University. His research focuses on data science and social computing, especially in the contexts of social/business networks and social media. His research has been covered by the *BBC*, *Washington Post*, *Forbes*, and others in more than twenty countries.

ABSTRACT: We focus on predicting the profitability of a movie to support movie-investment decisions at early stages of film production. By leveraging data from various sources, and using social network analysis and text mining techniques, the proposed system extracts several types of features, including “who” is in the cast, “what” a movie is about, “when” a movie will be released, as well as “hybrid” features. Experiment results showed that the system outperforms benchmark methods by a large margin. Novel features we proposed made weighty contributions to the prediction. In addition to designing a decision support system with practical utility, we also analyzed key factors of movie profitability. Furthermore, we demonstrated the prescriptive value of our system by illustrating how it can be used to recommend a set of profit-maximizing cast members. This research highlights the power of predictive and prescriptive data analytics in information systems to aid business decisions.

KEY WORDS AND PHRASES: decision support, movie investments, movie profitability, predictive analytics, prescriptive analytics, social network analysis, text mining.

The motion picture industry is a multibillion-dollar business. In 2015, the United States and Canada saw total box office revenues topping \$11.1 billion [29]. Nevertheless, the financial success of a movie is largely uncertain, with “hits” and “flops” released almost every year. While researchers have undertaken the task of predicting movie success using various approaches, they have attempted to predict

box office revenues or theater admissions. However, from an investor's standpoint, one would want to be as assured as possible that his/her investment will ultimately lead to returns. For instance, *Evan Almighty* earned a high gross revenue of \$100 million, but cost \$175 million to produce, whereas *Super Troopers* cost \$3 million, but earned \$18.5 million. The latter is certainly more appealing from an investment standpoint. In fact, among movies produced between 2000 and 2010 in the United States, only 36 percent had box office revenues higher than their production budgets, which further highlights the importance of making the right investment decisions. Therefore, our work defines a movie's success as its profitability, and attempts to predict such success in an automated way to better support movie investors' decisions.

The production process of a movie begins with the development phase, including the construction of a script and screenplay. Next, the potential film enters the preproduction phase, the most crucial to success. During this phase, the filmmaking team is assembled, filming locations are determined, and investments are obtained, among other decisions. Then, the film moves into the actual production phase, in which filming occurs. The postproduction phase involves the insertion of aftereffects and editing. The last phase is distribution [16]. To support the investment decisions of a movie, the prediction of profitability has to be provided before the actual production phase. In this research, we are interested in predicting a movie's financial success during its preproduction phase. Consequently, we can only leverage data that are available at this time.

Predictions made right before [17] or after [7, 26, 40] the official release (the final phase in movie construction) may have more data to use and get more accurate results, but they are too late for investors to make any meaningful decision. Building on previous work [22], this research proposes a movie investor assurance system (MIAS) to provide early predictions of movie profitability. Based on historical data, the system automatically extracts important characteristics for each movie, including "who" will be involved in the movie, "what" the movie is about, "when" the movie will be released, and the match between these features. It then uses various machine-learning methods to predict the success of the movie with different criteria for profitability.

The overarching research question of this study is to predict movie profitability using data available only during the preproduction stage of movie development. By proposing the first system to predict movie profitability at an early stage, the main contributions of this research are in two areas: first, this work demonstrates how freely available data of different types (including structured data, network data, and unstructured data) can be collected, fused, and analyzed to train machine-learning algorithms. When designing and developing information system artifacts [20, 31], such data-based approaches can provide powerful forecasts and recommendations to aid business decisions. To the best of our knowledge, we are also the first to leverage such data and models to prescribe profit-maximizing casts. Second, our research proposes several novel features, such as dynamic network features, plot topic distributions, the match between "what" and "who," the match between "what"

and “when,” and the use of profit-based star power measures to predict the profitability of movies at early stages. We showed that these features all make great contributions to the system’s performance, and help to explain important factors behind movies’ profitability.

Related Work

The Definition of Success

The way in which success is defined is of paramount importance to the problem, but past works have focused primarily on gross box office revenue [3, 4, 18, 28, 30, 34], while some have used the number of admissions [5, 26]. The basic assumption for using the two as success metrics is simple—a movie that sells well at the box office is considered a success. However, the two metrics ignore how much it costs to produce a movie. In fact, our analysis of historical data also found that revenues are not directly related to profits (more details will be presented in the Discussion section). Thus a more meaningful measure of success should be profitability, whether it is the numeric value of profits [36] or the return on investment (ROI) [14].

After a success metric was chosen, many studies categorized movies into two classes based on revenues (success or not) and adopted binary classifications as their predictive task; some considered the prediction as a multiclass classification problem and tried to classify movies into several discrete categories [30]. Meanwhile, predictions are also made on continuous numerical values of success metrics [17, 28, 38], with values of these metrics being logarithmized in several studies [34, 36, 40].

Features of Movie Success

The accuracy of a predictive model depends a lot on the extraction and engineering of features (aka, independent variables). When it comes to studying movie success, three types of features have been explored: audience-based, release-based, and movie-based features.

Audience-based features are about potential audiences’ reception of a movie. The more optimistic, positive, or excited audiences are about a movie, the more likely it is to have higher revenues, and vice versa. Movie reception can be retrieved from different types of media, such as Twitter [4], trailer comments [3], blogs [18], news articles [40], and movie reviews [26].

Release-based features focus on the availability of a movie and the time of its release. One such feature that captures availability at release is the number of theaters a movie opens in [28, 30, 32, 34, 38, 40]. The more theaters that show a movie, the more likely that the movie will have higher revenues. Many movies are targeted for release at a certain time. For example, holiday releases as well as seasons and dates of releases (spring, summer, etc.), are commonly used in the

prediction problem [8, 18, 30, 34]. Some studies also attempted to capture the competition at the time of release [18, 30], which could negatively affect revenues.

Movie-based features are those that are directly related to a movie itself, including who is in the cast and what the movie is about. The most popular feature for cast members is a movie's star power—whether the movie casts star actors. The star power of actors has been captured by actor earnings [30], past award nominations [7], actor rankings [34], and the actor's number of Twitter followers [3]. It was agreed that higher star power is helpful for a movie's success. However, no research has explored the profitability of actors. As it costs a great amount of money to cast a famous actor, we believe an actor's record of profitability will be a better indicator of a movie's profitability than the actor's record in generating revenues. Moreover, the role of directors in a movie's financial success is often overlooked or downplayed. While some research has investigated the individual success of directors [24], few studies have actually tried to connect directors' star power to movies' financial success. Some past studies have argued that the economic performance of movies is not affected by the presence of star directors [7], and directors' value is not as important as actors' for movie revenues [27]. Contrary to these selected past studies, we believe that both actors and directors are crucial for a film's success. Because directors, in particular, play important roles in movie production [24], our research will examine the effect of directors on movie profitability, in addition to actors.

In addition to individual actors and directors, the cast of a movie has also been explored from a teamwork perspective—whether individuals in a team can work together and develop “team chemistry” [26]. Studies of organizations and teams have revealed that team members' prior experience or expertise is beneficial for team success, while the diversity of a team helps too, especially in the context of bringing creative ideas and unique experience to teams for scientific research and performing arts [19, 35]. The diversity and the familiarity of a cast contribute to a director's success in receiving awards [24], and the movie's box office revenues [26]. Cast members' previous experience also positively influences revenues [27]. Nevertheless, there are several important limitations to consider. On one hand, many of the measures of teamwork were simplistic and problematic. For example, an actor's experience was based solely on the number of previous movie appearances, without considering what types of movies he or she has contributed to, and thus has more experience in. Also, team members' degree dispersions were used to reflect a team's diversity even though a team composed of actors who have never collaborated with each other can still feature a uniform degree distribution. Although the existence of structural holes can reflect a team's diversity, the measurement of structural holes was simplified to the density of a network. The two concepts are only very loosely related, however.

On the other hand, the data size was small in many studies. For instance, the top ten movies (by revenue) in each year (a total sample size of 160–180 movies) were studied in Meiseberg and Ehrmann [26] and Meiseberg et al. [27]. With such a small

sample, an actor's experience and previous collaborations cannot be completely captured. The selection bias toward more successful movies also hurt the validity of the results. Thus, in this research, we leveraged much larger data sets, derived new and more accurate ways to capture individual actor experience and team diversity, and related them to movie profitability.

In terms of what a movie is about, features such as genre, MPAA rating, whether or not a movie is a sequel, and run time have often been incorporated into success predictions, as well as in other domains [1]. Besides such meta data about a movie, to get a better idea of a movie's content, one needs to examine its plot or script. Two earlier studies leveraged the texts of movie scripts for success prediction [15, 16]. Some of the basic text-based features are easy to obtain, such as the number of words, and the number of sentences. However, more informative textual features in these studies depend on manual annotations by human experts, such as the degree to which the story or hero is logical, and whether or not the story has a believable ending. Because movie scripts can be very long, the manual annotations are time-consuming. Also, only a small number of movies' scripts are available in a uniform and professional format. Thus a predictive model based on features from scripts can only be trained on a small pool of movies, which may limit the predictive power for future movies. Thus an automated way to analyze openly available text-based movie content is necessary for a decision support system to learn from large-scale data sets.

For our research question of predicting movie profitability at an early stage, we cannot take advantage of most audience-based features and some of the release-based features because these are not available when making investment decisions. For instance, YouTube comments only appear after a movie trailer is released; likewise, the number of theaters a movie will be released in is not known until the end of the movie's production. In addition, these features from different groups were treated as standalone and independent, whereas the interaction or match between features from different groups, such as actors' star power along with their experience in different movie genres, or the popularity of a certain type of movie during a specific time period, can provide valuable information about a movie's success.

Therefore, we focus mainly on four types of features: "Who" features—who is involved in a movie, "When" features—when a movie will be released, "What" features from both meta data and text of movie plot synopses (movie plot synopses are openly available from most movie data archives, yet they can still reflect movies' content in the absence of a full script), and "Hybrid" features—the match between "What" and "Who" and the match between "What" and "When." Our feature set includes popular features from the literature (e.g., measuring actors' star power using their total gross revenues), new features proposed to better measure previously proven factors of movie success (e.g., team expertise and diversity), as well as features representing new factors that may be related to movie success (e.g., actor-director collaboration, and market trend by genre). All the features adopted by our system can be extracted in an automated fashion by using text mining and social network analysis techniques. In addition, from a theoretical perspective, this study also examined whether previous findings about star power of actors and directors

and about teamwork are still valid when movie success is measured by profit, instead of revenues, based on a much larger data set.

The System Framework

Figure 1 illustrates the framework of our MIAS. The first phase is data acquisition, because we based our prediction on historical data. We picked two popular and complementary sources—IMDb and Box Office Mojo. IMDb has better coverage of movie plot synopses, and Box Office Mojo, as its name suggests, provides more comprehensive data of movie revenues and budgets. In other words, the two data sources can be used jointly to acquire data for many movies. As for data collection methods, the two sources are different as well. IMDb has an application program interface (API) to provide movie data. The data from Box Office Mojo can only be obtained by the public from its web pages. To get a more comprehensive data set, our system employs two scripts: one interacts with APIs, while the other is a web scraper to retrieve and parse HTML data from web pages. We believe these two methods should be able to handle data from most open archives on the Internet.

In the second phase, data from both sources are cleaned, transformed, consolidated, and stored in a database. During this undertaking, we make sure that acquired data are put in a consistent format, and that the data are not duplicated within the database. For example, for movie titles, characters such as “*” and “-” are removed. Such standardization ensures that extraneous characters do not occlude the matching of titles between the two data sources. For plot synopses, the Porter stemmer was used, and stop words (such as “the”) were also removed.

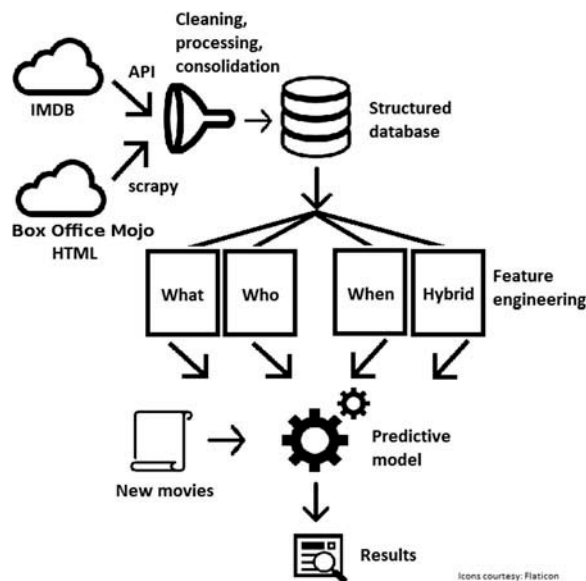


Figure 1. The Framework of MIAS

The third phase, “feature engineering,” involves using the acquired data to construct features that will ultimately be used to train a predictive model. Categorically, we classify various features into one of four groups—“what,” “who,” “when,” or “hybrid.” These feature groups will be described in detail later in the study. Features used in this study, and the reasons for including them, will be discussed at length in the next section.

With a reasonable and well-rounded set of features in place, a predictive model can be trained in the fourth phase of MIAS. Users of MIAS can define their own profitability metric or threshold based on the goals they have for their movies. They can employ cross-validation to select the best-performing model of prediction, along with its parameters, based on a maximized performance criteria of interest, such as overall accuracy, precision, or recall. The Experiments section will discuss our experiments in detail.

Feature Engineering

Based on historical data acquired from online archives, we derived four groups of features: “who” features, “what” features, “when” features, and “hybrid” features that match “who” with “what,” and “what” with “when.”

“Who” Features

Star Power

The very nature of the movie industry is characterized by people who make movies. Successful actors and directors are crowd favorites who are well-known throughout the world. Talented individuals can leverage not only their refined industry skills but also the associated “name brand effect,” which draws crowds and increases sales [4, 14, 37]. This effect is typically referred to as “star power.” Because our goal is to predict profitability, our star power features for a movie are based on its cast members’ records in generating both box office revenues and profits.

Tenure of an actor reflects how much experience he or she may have in the industry. It is calculated as the time difference (in years) between the movie in which an actor most recently appeared and that in which he or she first appeared. For each movie, we calculate the *average* and *total* tenure for its first-billed actors.

Actor Gross is how much revenue an actor has generated during his or her tenure. Each individual’s total gross is the sum of revenues from all the movies the actor has starred in, while an individual’s average gross is the actor’s total gross divided by the number of movies starred in. For each movie, we calculated the sum and average of total gross, as well as the average of actors’ average gross, for all first-billed cast members.

Director Gross measures the past success of directors. We calculated for each director the *total* and *average* gross for movies he or she has directed.

Actor/Director Profit measures the amount of profit an actor/director has earned through his or her career before the movie to be predicted. For each actor/director, we derived total profit, average profit, and top profit—the profit of the most profitable movie for the actor/director.

Network-based Features

Star power features listed above reflect whether or not a movie’s cast consists of senior and successful individuals (actors and the director). To capture team characteristics, we explored the avenue of social networks, which have the potential to yield a wealth of information about interpersonal interactions and collaboration [25, 41, 42], including teams for movie productions [23, 25].

For our predictive model, we constructed a dynamic collaboration network among actors based on their coappearances (i.e., costarring) in previous movies. In such a network, a node represents an actor. For any arbitrary year, an undirected edge was drawn between two actors if they costarred in a movie during that year. If an edge already existed between the two, indicating that they had collaborated in the past, the edge weight was incremented by 1. Therefore, the aggregated network for a given year includes all earlier years of collaborations, plus those that happen in that year. Figure 2 shows an example network.

Our network features consist of static features and dynamic features. When analyzing the team T_m for movie m in year y , the social network among the movie’s cast members were used to extract the following static features:

Network Heterogeneity: for each movie, we measured its team diversity by examining the structural network similarity between cast members. Specifically, based on each actor’s neighborhood vector in the adjacency matrix, we calculated the average cosine similarity between each pair of actors in the movie, which is denoted by Equation (1). In this equation let $|T_m|$ denote the number of cast members in team T_m for movie m ; $Act_i \cdot Act_j$ is the dot product between two actors on the team; and $\|Act_i\| \|Act_j\|$ is the magnitude of the two actor vectors. Higher similarity means team members have been working with similar peers (including one another), and vice versa. We believe this measure can capture previous collaborations among team members better than degree dispersion [28], which does not consider who an actor is connected to in a network.

$$H_m = \frac{1}{(|T_m|(|T_m| - 1)/2)} \sum_{i=1}^{|T_m|-1} \sum_{j=i+1}^{|T_m|} \frac{Act_i \bullet Act_j}{\|Act_i\| \|Act_j\|} \quad (1)$$

Average Degree represents the average number of unique collaborations for each cast member in a given movie. This metric is meant to capture the “degree” to which the team is truly bringing rich expertise and experience to the production of a movie [27].

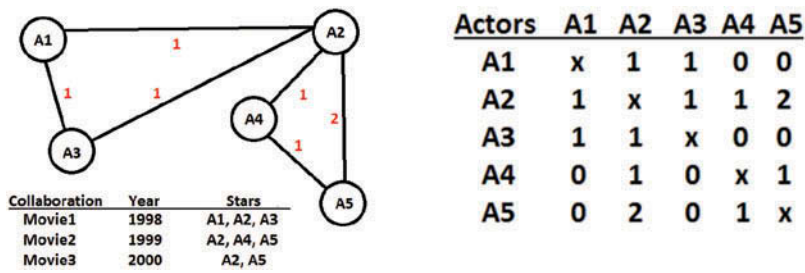


Figure 2. A Collaboration Network Example with the Network Structure (Left) and the Corresponding Adjacency Matrix (Right)

Total and Average Betweenness Centrality: In addition to those with many unique collaborators, “brokers” who can bridge different and otherwise less interconnected groups are also in a good position to bring in unique expertise and experience. These “brokers” often have high betweenness centralities and are said to have high social capital [10]. Having such “brokers” on a team can increase the team’s diversity by creating new ideas and producing innovations [9].

In addition to collaboration among actors, we also considered the collaboration between actors and directors by examining whether an actor and a director have worked together before and whether previous collaborations between them were successful.

Average Actor–Director Collaboration Frequency and Profitability of movie m is the average number of times that cast members of m have previously appeared in movies directed by the director of m , and the average profit per movie earned from all past collaborations between the actors and the director of movie m , respectively.

The features introduced above are based on the static structure of the year $y - 1$ network before movie m was produced in year y . Once movie m was produced in year y , its cast members formed a new team, which would add edges to the collaboration network and change its structure. Our newly developed dynamic network features tries to capture the spanning of structural holes after a new movie is produced.

Structural holes are an important concept in network analysis [10], and networks pertaining to movies are no exception [39]. Some studies suggest that a movie, which establishes interactor links that span structural holes, is more likely to succeed [7], although they do not quantify the degree to which a movie spans structural holes. Some research has used the clustering coefficient for each team member to measure the existence of structural holes [27]. The clustering coefficient of a node is the probability that the node’s neighbors are also connected to each other. It measures how a node’s neighborhood is “clustered” together. A network clustering coefficient is the average of all its node values of this metric. However, the static value of a team’s average clustering coefficient alone, in the current collaboration network, can only show structural holes at the ego level (i.e., among immediate neighbors of an ego). To capture the spanning of structural holes at the network

level, we used the following two dynamic network features to measure how the structure of the collaboration network changed after incorporating the new collaboration.

Decrease in clustering coefficient: A new movie will add edges to an existing collaboration network. If these new edges connect nodes that are originally only two hops away from each other, then the clustering coefficient of the network will increase, but such edges only reinforce existing clusters. However, by creating new neighbors without closing the triad, new edges that connect nodes that are originally far away will decrease the network’s clustering coefficient. It has been found that decreasing the clustering coefficient of a social network can facilitate the diffusion of information across the network by breaking up existing clusters [41]. Thus we included the decrease in the clustering coefficient of a collaboration network after forming the team for movie m to measure whether new collaborations can break existing clusters.

Decrease in average shortest path: We also proposed to use how the production of a movie m decreases the average shortest path length of the social network, because adding edges that span structural holes usually significantly decreases such path length. Specifically, after adding to $Network_{y-1}$ edges that correspond to the cast of T_m produced at year y , we calculated how much the average shortest path length of the new network decreased, compared to $Network_{y-1}$. The more such length decreases, the more movie m ’s cast can span structural holes in $Network_{y-1}$.

“What” Features

In addition to “who” is in the cast, another natural and important indicator of a movie’s future profitability is what the movie is about. Such information is usually available with high certainty prior to movie funding efforts. To reflect what a movie is about, the “what” features in our model include meta features, such as *genre* (e.g., action, sci-fi, family) and *rating* (e.g., PG13, R, etc.), represented as binary categorical variables.

We also included a fine-grained description of a movie’s content—its plot synopsis. While the full script of a movie is a better representation of the movie’s content, such scripts are very difficult to obtain for a large number of movies, especially those that were not very successful. Thus we used plot synopses as approximations for full scripts, in that plot synopses are usually publicly available. This allows our predictions to be based on a larger pool of movies.

Representing texts from plot synopses with traditional unigrams and bigrams will have high dimensionality and suffer from sparsity. At a higher level, topic modeling techniques, such as latent Dirichlet allocation (LDA) [6], can give a better picture of what a plot is about. The input for LDA is a textual corpus of plot synopses and the output is a group of topics, each being represented by a probabilistic distribution over archetypal words. Those words, having a high probability of a given topic, are considered representative keywords for that particular topic. Each plot synopsis is

also assigned a probabilistic distribution over all the topics. In the topic distribution vector of a movie's plot, each element is the probability that the movie represents each topic. Therefore, each element is a continuous numerical value $\in [0, 1]$, where 0 indicates that the movie does not at all represent the topic, and 1 indicates a perfect representation. This topic distribution reflects the content of the movie at an aggregated level and can be used as features for predictions.

In addition to these topics derived from LDA, some movies' plots are adaptations from other sources, an important consideration, especially when the original source has achieved certain levels of success. For example, *The Hunger Games* and *Harry Potter* are both adapted from best-selling novels. As such, one of our "what" features was about adaptations: whether a movie's plot was adapted from a comic, a true story, or a book/novel.

"When" Features

With the movie industry being an avenue for entertainment, its market sees peaks and declines over time, which can suggest how well an already produced movie might fare in the future. Thus we incorporated the following "when" features in our model: *Average Annual Profit* is the average profit across all movies in the year prior to the planned release of movie m . It captures the overall profitability of the movie industry before a movie is released. *Release dates* combines several features of when a movie will be released, including whether it will be a holiday release and which season of the year (spring, summer, fall, winter). While a holiday or summer release may attract a larger audience and thus generate more revenues [2], it also requires a larger budget for marketing and distribution during these competitive periods. Although the exact release date is not completely definitive before filming, a target trajectory usually exists at preproduction stages.

Hybrid Features

Besides standalone features of "who" is in a movie, "what" a movie is about, and "when" a movie will be released, it is also important to capture the "match" between these features. Our hybrid features try to reflect such matches between "what" and "who" as well as between "what" and "when." For example, it may be important to form a team of actors based on their previous experience with the genre of the movie being planned instead of just their star power. Similarly, the investment in a movie whose genre is gaining popularity may increase the chance of success.

"What" + "Who"

In observing the movie industry and the actors, we can distinguish various so-called roles that these actors seem to adopt. For instance, Seth Rogan is typified

by his appearance in comedies, and Arnold Schwarzenegger exhibits proficiency as an action movie star. Should a movie then try to include those who have extensive experience in its genre? Or conversely, does a surprising cast draw a greater audience to theaters (e.g., having Schwarzenegger in a comedy or a romance)? Although these questions have not been addressed in the literature, we believe that better measurements of an actor's expertise with regard to genres can help us more accurately determine the expertise and diversity of a movie's cast.

To measure an actor's previous experience and expertise in movies of different genres we define, for each actor j , a genre experience vector $A_j = [a_{j,1}, \dots, a_{j,k}, \dots, a_{j,K}]$, where $a_{j,k}$ is the proportion of the number of times actor j appeared in movies of genre k . A total of $K = 26$ unique genres are defined. Similarly, a movie m is also represented as a genre vector $G_m = [g_{m,1}, \dots, g_{m,k}, \dots, g_{m,K}]$, where $g_{m,k} = 1$ indicates that movie m has genre k , and $g_{m,k} = 0$ otherwise. Note that some movies can have more than one genre. For example, the genre of *Spiderman* is both action and adventure. By measuring the similarity between actors' genre experience vectors and movies' genre vectors, we designed several features that address the genre-based expertise brought by cast members to a given film m 's team T_m .

Average Genre Expertise (AGE): captures the average cast experience with respect to the current movie's genre. Movie m 's *AGE* is defined in Equation (2).

$$AGE_m = \frac{1}{|T_m|} \sum_{j=1}^{|T_m|} G_m \bullet A_j \quad (2)$$

Weighted Average Genre Expertise (WAGE) extends *AGE* by incorporating an actor's star power, measured by actor gross, in each genre. As defined in Equation (3), the *WAGE* of movie m is essentially the movie's *AGE* weighted by each cast member j 's gross revenue R_j . In other words, a movie with a big star who is familiar with its genre will have high *WAGE*.

$$WAGE_m = \frac{1}{|T_m|} \sum_{j=1}^{|T_m|} \log(R_j) * (G_m \bullet A_j) \quad (3)$$

Cast Novelty is defined in a way similar to *WAGE*. While *WAGE* is an average value that tries to capture a cast's experience in the movie's genre, cast novelty focuses on team diversity—whether a movie has a big star who has rarely appeared in movies of this genre before. It is the maximum value among all actors' star-power-weighted inverse experience in movie m 's genre (Equation [4]). Higher values indicate having an unexpected star appear in a given movie.

$$CN_m = \max \left\{ \frac{\log(R_j)}{G_m \bullet A_j + 1}, \forall j \in T_m \right\} \quad (4)$$

“What” + “When”

Similar to the overall market volume for movies, which changes over time, consumers’ preferences of movies may also evolve from year to year. For example, while movies like *American Pie* and *National Lampoon’s Van Wilder* were popular in the late 1990s and early 2000s, moviegoers have recently been flocking to horror movies, such as *Paranormal Activity*, and those characterized by superheroes, such as *The Avengers*. Although the latter category is nothing definitively new to the silver screen, the movie industry has seen greater levels of success in recent years with this particular focus and, as such, a greater influx of such movies. Meanwhile, competition may also affect the profitability of movie m because other movies released during a similar time period may detract from movie m ’s viewer base [27]. Thus, in addition to capturing “when” a movie will be released, we also consider movies of a similar genre performed in the previous year, as well as the level of competition during a movie’s planned release time.

Annual Profitability Percentage by Genre is the percentage of movies of the same genre as movie m , in the year prior to the planned release of movie m , that were profitable. This feature reflects the degree of success for movies that share the same genre as the movie being considered.

Annual Weighted Profitability by Genre (AWPG) is derived from movie genre vectors defined earlier in this study. For movie m in year y , the profitability of each movie m' in year $y - 1$ are summed up and weighted by the cosine similarities between genre vectors of m and each m' . Equation (5) illustrates how to calculate the AWPG for movie m in year y , where G_m is the genre vector for movie m and $p(m')$ is the profitability of movie m' . This feature indicates the overall previous-year profitability of movies whose genre is similar to a given movie.

$$AWPG_m = \sum_{m' \in y-1} sim(G_m, G_{m'}) * p(m') \quad (5)$$

Competition reflects the other movies that will be released during a similar time period. It is calculated by considering the average star power of all other movies released within one month of movie m ’s release date. This feature indicates the degree to which other big-name stars are appearing in movies at a similar time, which may detract from movie m ’s viewership. The inclusion of such a feature is based on the fact that, even prior to production, a movie has a loose, or at times even definitive, trajectory for release. By defining competition to be within one month (plus or minus) of the original release date, such a notion of “approximate release” is maintained, even in instances in which the exact date may have been well-established (i.e., such an approach is conservative).

Experiments

Data Set and Basic Statistics

Our original data set, collected from both Box Office Mojo and IMDb, consisted of 14,097 movies, along with 4,420 actors. While movies in our data set date back to 1921, we focused our study on movies released during the eleven-year period of 2000 through 2010 because this period is recent enough to reflect the current state of the industry, while ensuring that a sufficient amount of time has elapsed since the movies were released so that revenue data can be accurately updated.

As our goal is to predict movie success measured by profits, our data set for experiments included only those movies that have both budget and box office revenue data available. We also excluded movies of an “Unknown” genre, or an “Unknown” MPAA rating. “Documentary” genre movies were also excluded because they are typically not released to theaters and may not involve professional actors. In addition, any movie designated as being part of a franchise, a sequel, or a remake was also excluded (e.g., *Iron Man*, *Iron Man 2*, etc.). We made this decision because the success of a sequel can depend heavily on the success of earlier movies in the same franchise. Also, the content of sequels and remakes and their selection of cast members are also highly limited by earlier counterparts. Thus what is behind the success of a sequel or remake may be very different from that of other movies.

With these considerations in mind, our final data set for experiments consisted of 2,506 movies. A distribution by genre of these 2,506 movies, relative to all movies released during the period, is presented in Figure 3. The distribution suggests that our data set is a representative sample overall, with the exception of the “Foreign” genre. This makes sense because budget and revenue data may be more difficult to obtain for movies that are produced, and in all likelihood, released outside the United States. Based on the plot synopses of these movies, we used LDA to generate thirty topics. Top keywords of these topics are listed in Table 1.

While the experiment will predict the success of 2,506 movies during an eleven-year period, the collaboration network we built for this study incorporates the collaboration between all actors in all 14,097 movies in our data set. The initial unweighted, undirected network was aggregated to the year 1999, with networks for

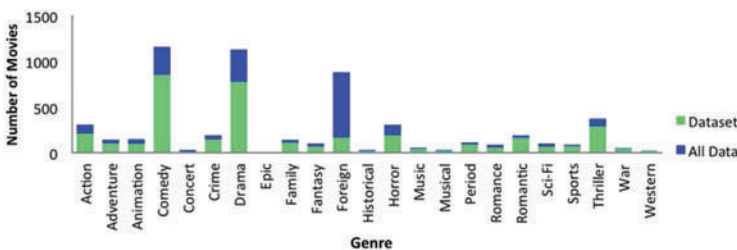


Figure 3. Distribution of Movies by Genre (2000–2010)

Table 1. Topics and Keywords Generated by LDA from Plot Synopses

| Topic | Keywords | Topic | Keywords |
|-------|---|-------|---|
| 1 | Wife, husband, marriage, child, couple | 16 | Man, young, become, past, truth |
| 2 | People, movie, story, show, tv | 17 | He, want, she, know, tell |
| 3 | Back, even, good, time, start | 18 | War, mission, American, government, fight |
| 4 | Music, band, famous, star, place | 19 | Love, young, women, heart, marry |
| 5 | First, world, people, state, country | 20 | Team, game, win, dream, big |
| 6 | Money, back, plan, help, deal | 21 | Work, job, business, company, success |
| 7 | Man, begin, believes, situation, hospital | 22 | School, high, parents, boy, girl |
| 8 | Life, young, city, world, lives | 23 | Friend, girlfriend, party, boyfriend, college |
| 9 | Find, way, help, search, journey | 24 | Story, film, based, documentary, history |
| 10 | Group, find, survive, crew, remote | 25 | Police, murder, drug, prison, kill |
| 11 | One, life, never, day, always | 26 | Two, lives, relationship, together, sex |
| 12 | World, stop, evil, power, battle | 27 | He, find, she, finally, arrive |
| 13 | Family, father, son, mother, home | 28 | Years, time, later, death, since |
| 14 | Night, day, car, trip, train | 29 | Events, forced, act, unexpected, secrets |
| 15 | New, life, deam, everything, lost | 30 | Town, local, small, gang, store |

subsequent years being updated to reflect that year's new collaborations. In all, we created eleven snapshots of the collaboration network from 1999 to 2009.

The Measure of Success

To predict whether a movie is successful, we measured a movie's profitability using two metrics—raw profit and return on investment (ROI). Raw profit is simply revenue minus budget. For movies in our data set, only 36 percent of the movies had positive profits. However, a profit of \$10,000 from a movie that costs \$1 million to produce is certainly not an attractive investment prospect. Thus in our experiments, we also adopted return on investment (ROI) as in [14]. Considering both profit and budget, ROI is defined as $ROI = \frac{Revenue - budget}{Budget}$. The higher the ROI is, the more profitable a movie is, and vice versa. To avoid “trivial” and disinteresting profit returns, we raised the profit bar for a movie to be considered “truly profitable,” as we will discuss in our experiments.

Interestingly, the data suggest that profitability, as measured by ROI, is not necessarily reflected by box office revenues. The correlation coefficient between revenues and ROIs is only 0.077. In other words, having a great box office revenue does not necessarily mean a high ROI. Similarly, albeit with a higher correlation

coefficient, raw profit has a correlation coefficient of 0.22 with revenue. These factors further highlight the need for an accurate prediction of profitability.

Classification of Profitability

The prediction of a movie's success can be modeled as a classification problem to decide whether a movie should be considered a success or not based on either ROI or profit. Although there is no agreed-upon industry "gold standard" of an ideal ROI or profit, it is reasonable to assume that one would like to see some substantive returns from a successful movie, given that millions of dollars are invested with considerable risks. Also, any form of profit is better than a loss. Thus we elected to define the decision boundary between successful and unsuccessful movies in two different ways for both binary and multiclass predictions.

For both binary and multiclass classification of movie success, we tried a variety of algorithms, including logistic regression, naive Bayes, support vector machines (SVM), multilayer perceptron (MLP), decision trees (J48), random forest, and the LogitBoost, and selected the one with the best overall performance based on the following four metrics (using ten-fold cross-validation), where higher values indicate better performance: (1) *Classification accuracy*, which is the percentage of correctly predicted instances; (2) *Precision* (positive class), which is the number of instances classified as being positive that are actually successful, divided by the number of instances classified as being successful; (3) *Recall* (positive class), which is the number of instances classified as being positive that are actually successful, divided by the number of instances that are actually successful; and (4) *The Area Under the Receiver Operator Characteristic Curve (AUC)*. The curve plots the true positive rate against the false positive rate. An *AUC* of 1 means a perfect classification whereas 0.5 refers to a random guess. Being more robust against prior distributions, *AUC* is considered by many to be one of the best indicators of a classifier's performance [33]. In multiclass classifications, we reported a weighted average of *AUCs*: first calculate the *AUC* obtained for each class, and then weight each *AUC* according to the number of instances that fall into each of these classes relative to the total number of instances.

In addition to identifying the best-performing algorithm, we also evaluated whether and how features we proposed in this research contributed to the prediction. We included, in a "New" feature group, those novel features that we proposed and are used for the first time to predict movie success. Features in this "New" group include (1) features related to actor and director profits, actor-director collaboration, and dynamic network features (e.g., decrease in the average shortest path length) from the "Who" group; (2) topic distribution features from the "What" group; (3) average annual profit in the "When" group; and (4) all features in the "Hybrid" group.

To further evaluate the performance of our predictive model, we also compared it with two benchmark models. Our benchmark models were constructed using

features that were defined in past studies and derived from our data. We also used our definitions of profitability and the same set of classification algorithms. In so doing, we were able to estimate how these past studies would have performed, allowing for an apples-to-apples comparison. Benchmark 1 was based on Vany and Walls [36] and Benchmark 2 was based on Walls [38]. Among previous studies of box office revenue predictions, we selected these two because most features used in their studies are available prior to a movie's production, which is similar to our early prediction problem. For example, the following features were used in Vany and Walls [36]: star power, sequel, genre, rating, and year of release. Similarly, Walls [38] used film budget/cost, number of screens the film was released on, sequel, star power, genre, and rating. We excluded the sequel feature, as our data set excludes such movies, and the number of screens feature, because this information is not available prior to release. To make the comparison consistent, we used our definition of star power for the two benchmark models. We reported results of best-performing classifiers for both benchmark methods, along with the performance of our approach.

We would like to point out, prior to reporting the results of such experiments, that a user of the MIAS system has the option to select the desired algorithm used in making predictions based on any of the aforementioned criteria they are interested in maximizing. Investors are most likely interested in recall and *AUC* based on their preferences. Selecting an algorithm that maximizes recall helps ensure that investors will not miss out on any movies that end up being profitable, but may end up with movies that are not eventually profitable. An algorithm that maximizes *AUC*, on the other hand, trades off ensuring that investment opportunities are not missed, but has the added benefit of protecting the investor from a monetary loss.

Furthermore, we imagine that users of this system will be tailoring, tuning, and evaluating its performance as new data are incorporated. In so doing, the user of MIAS may find that another algorithm performs better in terms of a particular criterion of interest. If this is the case, the user, of course, has the option of making decisions based on the better-performing model.

Binary Classification

For binary classifications, a movie is classified into one of two classes: successful or unsuccessful. Two decision boundaries are evaluated and both ensure that a sufficient amount of either ROI or raw profit is obtained for a movie to be considered successful. The first decision boundary entails that a movie is considered successful if its ROI is within the top 30 percent of all movies. For our data set, this threshold translates to $ROI \geq 24$ percent. Table 2 lists the performance of the top two classification algorithms, a random forest classifier ($n = 200$) and a LogitBoost classifier.

Table 2 also highlights the contribution of "New" features to the prediction. When "New" features were removed, *AUC* and accuracy of the classifiers deteriorate 29

Table 2. Top Two Prediction Results of Our Binary Classification Model and the Performance Without “New” Features (with top 30 percent ROIs as the decision boundary)

| Classifier Model | Random forest | | Logit Boost | |
|---------------------|---------------|------------------|-------------|------------------|
| | Full model | W/o New feauters | Full model | W/o New Feauters |
| AUC | 0.863 | 0.616 | 0.833 | 0.653 |
| Accuracy | 0.834 | 0.675 | 0.812 | 0.697 |
| Precision | 0.82 | 0.454 | 0.844 | 0.492 |
| Recall | 0.575 | 0.380 | 0.465 | 0.129 |

percent and 20 percent, respectively, for random forest, and 22 percent and 15 percent for LogitBoost. Precision and recall also drop greatly. In addition, the top two classifiers of the two benchmark models (Table 3) trail our model in all four performance metrics. For instance, *AUC*s of the two benchmark models are 19 percent and 25 percent, respectively, lower than ours.

The second boundary we tested is $Profit \geq \$7.3$ million, which corresponds to one-quarter standard deviation above the mean profit. With this threshold, 21.4 percent of the movies in our data set were considered successful. Top performing algorithms are able to reach *AUC* and accuracy over 0.9 (see Table 4). At the same time, our “New” features still make great contributions to the classification—the removal of these features dropped the *AUC* of the best-performing random forest classifier by 24 percent, and the LogitBoost classifier’s *AUC* decreased by 20 percent. A similar decrease can be found for accuracy, precision, and recall. Our model also retains an advantage over the two benchmark models (Table 5), leading by 22–27 percent *AUC*.

Multiclass Classification

In the case of multiclass classifications, we defined three possible classes for a movie: positive (“success”), negative (“failure”), or neutral (“average”) to provide

Table 3. Top Two Prediction Results for Benchmark Binary Classification Models (with top 30 percent ROIs as the decision boundary)

| Classifier | Benchmark 1 | | Benchmark 2 | |
|------------|---------------------|----------------|---------------------|-------------|
| | Logistic regression | Naive Bayesian | Logistic regression | Logit Boost |
| AUC | 0.672 | 0.651 | 0.701 | 0.651 |
| Accuracy | 0.702 | 0.686 | 0.724 | 0.686 |
| Precision | 0.516 | 0.475 | 0.603 | 0.475 |
| Recall | 0.188 | 0.367 | 0.252 | 0.367 |

Table 4. Top Two Prediction Results of Our Binary Classification Model and the Performance Without “New” Features (with $Profit \geq \$7.3$ million as the decision boundary)

| Classifier | Random forest | | Logit Boost | |
|------------|---------------|------------------|-------------|------------------|
| | Full model | W/o New features | Full model | W/o New features |
| AUC | 0.921 | 0.707 | 0.917 | 0.735 |
| Accuracy | 0.904 | 0.749 | 0.891 | 0.796 |
| Precision | 0.874 | 0.399 | 0.855 | 0.583 |
| Recall | 0.646 | 0.338 | 0.593 | 0.164 |

Table 5. Top Two Prediction Results for Benchmark Binary Classification Models (with $Profit \geq \$7.3$ million as the decision boundary)

| Classifier | Benchmark 1 | | Benchmark 2 | |
|------------|---------------------|-------------|---------------------|-------------|
| | Logistic regression | Logit Boost | Logistic regression | Logit Boost |
| AUC | 0.754 | 0.726 | 0.756 | 0.725 |
| Accuracy | 0.786 | 0.795 | 0.793 | 0.761 |
| Precision | 0.500 | 0.597 | 0.547 | 0.436 |
| Recall | 0.175 | 0.132 | 0.194 | 0.397 |

more information to investors on where they could expect a movie to fall with regard to profitability. For the multiclass prediction, we explored the imposition of cost [13] associated with misclassification, because the three classes are ordinal. In other words, not all misclassification errors are equally severe. For example, for investment decision support, predicting a failure to be a success would be worse than predicting it to be a neutral movie. The cost matrix for the multiclass classification is presented in Table 6—the penalty imposed for classifying a *successful* movie as a *failure* is 2, and vice versa, whereas the penalty for misclassifying by only one ordinal category (i.e., success as neutral, etc.) is 1.

Similar to binary classifications, we defined the three classes of success in two ways. The first way was to split movies into three equal-size classes: the positive class consists of movies with top one-third ROIs ($ROI \geq 10$ percent), the negative class consists of movies with the bottom one-third ROIs ($ROI \leq -78$ percent), and the other middle one-third into the neutral class ($-78 \text{ percent} < ROI < 10 \text{ percent}$). The second way was to classify movies with the top one-fourth ROIs as positive ($ROI \geq 47$ percent), the bottom one-fourth ROIs as negative ($ROI \leq -91$ percent), and the rest as neutral ($-91 \text{ percent} < ROI < 47 \text{ percent}$).

After comparing the performance of several classification models (including J48, Naive Bayes, MLP, SVM, logistic regression, and Logit Boost), random forest still emerged as the best classifier for both decision boundaries. Table 6 lists the

Table 6. Multiclass Classification Results of Our Model from the Best-Performing Random Forest Classifier, and the Performance Without “New” Features

| Measure Model | The 1st decision boundary | | The 2nd decision boundary | |
|--|---------------------------|------------------|---------------------------|------------------|
| | Full model | W/o New Features | Full model | W/o New Features |
| AUC | 0.847 | 0.636 | 0.85 | 0.657 |
| Accuracy | 0.679 | 0.459 | 0.73 | 0.508 |
| Precision (“Pos. Class”) | 0.769 | 0.483 | 0.803 | 0.435 |
| Recall (“Pos. Class”) | 0.711 | 0.482 | 0.671 | 0.424 |
| Total cost | 986 | 1882 | 732 | 1505 |
| “Pos. Class” refers to “Positive Class”. | | | | |

Table 7. Multiclass Classification Results of Benchmark Models Using Random Forest Classifiers

| Measure | The 1st decision boundary | | The 2nd decision boundary | |
|--------------------------|---------------------------|-------------|---------------------------|-------------|
| Model | Benchmark 1 | Benchmark 2 | Benchmark 1 | Benchmark 2 |
| AUC | 0.77 | 0.626 | 0.806 | 0.657 |
| Accuracy | 0.578 | 0.448 | 0.651 | 0.508 |
| Precision ("Pos. Class") | 0.474 | 0.456 | 0.473 | 0.406 |
| Recall ("Pos. Class") | 0.452 | 0.476 | 0.362 | 0.383 |
| Total cost | 1534 | 1915 | 1140 | 1509 |

"Pos. Class" refers to "Positive Class".

performance measures and Table 7 shows the performance from the two benchmark methods. Similar to that of binary classifications, our model outperforms the two benchmark models by reducing the total misclassification cost by 36–52 percent. Meanwhile, "New" features continue making great contributions to the prediction—the exclusion of these features from our model doubles the misclassification cost.

Discussion

Regression Analysis

While a random forest classifier can do a good job of predicting whether a movie will be successful, it is also important to understand the factors behind such success. A regression model will help us better assess the degree to which individual features influence predictive results, and to examine whether they are indicative of movie profitability. Furthermore, a regression model can also provide predictions on numeric values, in case the classification of movies into two or three discrete groups is not sufficient for investors' needs. Thus we also explored the prediction of continuous ROI values. It is worth noting that because the distribution of ROI is highly skewed, we applied a logarithm transformation to ROI, in the format of $\log(ROI + 1)$ as in Eliashberg et al. [16].

We tried six different algorithms, namely, LASSO, support vector regression (SVR), ridge regression, CART, M5P trees, and REP tree. Among them, we were particularly interested in LASSO and ridge regression for two reasons. First, coefficients of each feature in these models are able to offer valuable insights into how each feature contributes to a movie's profit. Second, multicollinearity may exist among our features (or independent variables). For example, the correlation between *Total Actor Profit* and *Average Actor Profit* is 0.96 (p -value < 0.001). LASSO and ridge regression both use regularization ($L1$ and $L2$, respectively) to penalize the nonzero values of the regression coefficients. Such regularization allows the model to select features that are more informative for the prediction and reduce the impact of collinearity [21].

Table 8 compares root mean squared errors (RMSE) of the six algorithms, with LASSO being the best at predicting numeric ROI values. Thus we used coefficients from LASSO to reveal factors behind movie success. To obtain these coefficients, we iteratively increased the penalizing λ value until all attribute-wise variance inflation factors were reduced to below 10 [12]. After achieving such a result with $\lambda = 0.0065$, 48 out of the 120 features in the LASSO model ended up having nonzero coefficients: 16 have negative coefficients, and 32 have positive coefficients.

Table 9 shows how many of the forty-eight features are from each feature group, including the “New” feature group for novel features proposed in this research. It turns out that the forty-eight features cover all feature groups, and more than half of them are “New” features. Besides the “New” feature group, the “What” group contributes the most features. For example, twelve of the thirty topics derived from LDA are among the forty-eight.

Table 10 lists the top five features by the value of their coefficients. A feature with a positive coefficient indicates that it has a positive influence on profitability, while a feature with a negative coefficient is indicative of a negative influence on profitability. For example, an increase of one standard deviation in the average profit of actor–director collaboration is associated with a 14.3 percent increase in $\log(ROI + 1)$, while being an R-rated movie decreases the $\log(ROI + 1)$ by 5.8 percent. While those with top negative coefficients are all “What” features, including genre (drama

Table 8. Results for Predicting $\log(ROI + 1)$ with Various Algorithms (ten-fold cross validation)

| Measure \ Algorithm | LASSO | SVR | Ridge regression | CART | M5P tree | REP tree |
|---------------------|-------|-------|------------------|-------|----------|----------|
| RMSE | 0.878 | 1.180 | 1.10 | 1.232 | 0.906 | 0.929 |

Table 9. Number of Features from Each Feature Group for the Forty-Eight Features with Nonzero Coefficients in LASSO

| Feature group | Number of features |
|----------------------|--------------------|
| Who (star power) | 7 |
| Who (Network-based) | 4 |
| What | 24 |
| When | 2 |
| Hybird (What + Who) | 9 |
| Hybird (What + When) | 2 |
| “New” features | 27 |

Table 10. Top Features with the Highest Positive and Lowest Negative Regression Coefficients from the LASSO Model

| | Feature group | feature | Coefficient |
|---------------------------|----------------------|--|-------------|
| Top positive coefficients | Who (star power) | Avg. profit of actor-director collaboration* | 0.143 |
| | Who (star power) | Avg. Director Gross | 0.039 |
| | When | Winter Release | 0.036 |
| | Who (star power) | Total Actor Profit* | 0.035 |
| | Hybird (What + When) | Annual Profit % by Genre* | 0.033 |
| Top negative coefficients | What | R rating | -0.058 |
| | What | Drama genre | -0.012 |
| | What | Topic 18* | -0.012 |
| | What | Topic 4* | -0.011 |
| | What | Foreign Genre | -0.009 |

and foreign), “R” rating, and plot topics related to wars and music, it is a mix of other feature groups (Who, When, and Hybrid) that are indicative of profit.

We also checked the robustness of the forty-eight top features by comparing them (generated by LASSO with $\lambda = 0.0065$) with top features generated by LASSO with two other λ values $\lambda = 0.0044$ and $\lambda = 0.0025$. With lower λ values in LASSO, more features are selected by the two models. However, all forty-eight features still appear in both lists of top features. These results confirmed that the list of forty-eight features is a robust representation of key factors for movie profitability.

Star Power and Movie Profits

As we mentioned before, star power features have a large and positive bearing on the success of movies. While previous studies agreed that higher start power is generally associated with movie success [36], they relied on movies’ box office revenue with actors’ star power measured by their total gross revenue. Figure 4 plots total actor gross against movie revenues and profits in our experiments. As we can observe, although total actor gross is moderately correlated with movie revenues (Pearson correlation coefficient 0.46), the correlation is much weaker with profits (Pearson correlation coefficient 0.16). In other words, having actors who have earned big box office revenues in a movie does not necessarily mean more profit for the movie. Such a result further highlights the difference between measuring movie success with revenues and profits.

By focusing on historical profitability records of both actors and directors, our results revealed some interesting findings about movie profitability. For example, the director is an important factor for movies’ profits. The top feature from our LASSO model is actually the average profit of previous actor–director collaboration. Also, the average profit of actor–director collaboration is a better indicator of $\log(ROI + 1)$

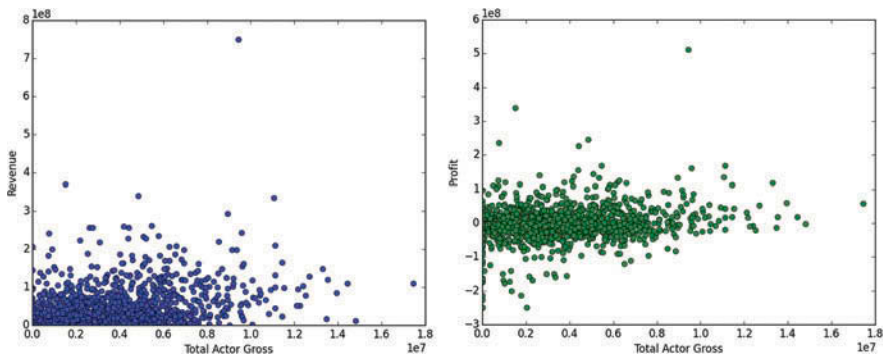


Figure 4. Total Actor Gross vs. Movie Revenues (Left) and Profits (Right)

than the traditional star power measure of total actor gross—the rank correlation between the average profit of actor–director collaboration and $\log(ROI + 1)$ is 0.47, while total actor gross has a rank correlation of 0.29 with $\log(ROI + 1)$. Furthermore, according to Table 10, having a star director is more indicative of profit than having a cast of star actors. Such findings actually contrasted with a few studies that even considered the effect of directors on movie success, albeit measured by box office revenues. We conjectured that the difference may be due to measuring movie success using profits instead of revenues, and the usage of a larger data set with movies whose success levels vary greatly. Although further investigations in this direction are beyond the scope of this research, we believe that this is an interesting result worth exploring from a team performance or marketing perspective.

When it comes to predicting movie profits, actors’ star power is better measured by the actors’ historical profits than with their gross. In fact, when we ranked actors by their total gross revenues and total profits, the ranking correlation is only moderate (Spearman coefficient of 0.60). Table 11 lists the top ten actors by total revenues and total profits, respectively, and there is only one (Julie Andrews) who appears on both lists. Although “big-name” movie stars are

Table 11. Top Ten Actors by Total Revenues and Total Profits

| Rank | By Total Revenue | By Total Profit |
|------|---------------------|-------------------|
| 1 | Clark Gregg | Orlando Bloom |
| 2 | Julie Andrews | Elijah Wood |
| 3 | Dakota Fanning | Robert Pattinson |
| 4 | Ashton Kutcher | Zoe Saldana |
| 5 | Steve Carell | Mike Myers |
| 6 | Morgan Freeman | Alan Rickman |
| 7 | Johnny Depp | Julie Andrews |
| 8 | Anna Kendrick | Samuel L. Jackson |
| 9 | Bryce Dallas Howard | Gary Oldman |
| 10 | Emma Roberts | Rupert Grint |

likely to attract a large audience, in the case of generating profits, the cost of casting such a star may not always be recouped via ticket sales.

Teamwork and Profitability

From a teamwork perspective, we investigated the effect of expertise and diversity, and found that both can contribute to a movie's profits. Our new metric to measure how much expertise a cast has in a specific movie's genre—*Average Genre Expertise*—is positively related to profits (with a coefficient of 0.007). The positive effect of Average Genre Expertise, along with the top positive coefficient for average actor-director collaboration profits, has highlighted the importance of a cast's expertise and successful collaboration experience in the past. At the same time, diversity is also a positive predictor of profits. Among our diversity metrics, decrease in clustering coefficient, cast novelty, and the spanning of structural holes (measured by the decrease in the average shortest path lengths) all have positive coefficients in the LASSO model. In other words, a cast with members who have previous experience in a movie's genre, yet with some fresh faces, is beneficial for a movie's profits.

Cast Selection via Inverse Classification

To demonstrate the power and additional utility of our proposed system beyond prediction, we extended our experiments to prescriptive analytics. Specifically, we show how the MIAS system can be leveraged to prescribe a set of cast members that maximize the probability of observing profitable returns according to one of our two definitions of profitability: $ROI \geq 24$ percent. Practically speaking, however, a thorough exploration of such prescriptive analysis is a massive undertaking, because the search space is large due to the large number of possible combinations of actors for a movie's cast; the optimization depends on the outcome of predictions as well. As such, we deployed a small-scale experiment in this study, leaving more definitive elaborations for future work.

The problem of selecting a profit-maximizing cast will be addressed through a process known as inverse classification [2, 11, 23]—the process of making perturbations to a test instance (in this case, a movie) such that the probability of some ideal class (i.e., positive profit) is maximized. The problem can be formulated as the partitioning of features into two groups: those that are unchangeable and those that are changeable [23]. For this experiment, we assume that cast member-based features are changeable (e.g., past actor-director collaborations, network-based features, etc.) and all others are unchangeable (e.g., genre, rating, etc.).

To construct an appropriate set of cast members, we used the original cast of a particular movie as the starting point, and observed retrospectively what changes to cast members could have led to higher probabilities of profitability. Specifically, for each of the original cast members i who planned to appear in movie M ($i \in C_M$), we wanted to construct a separate set of cast members to maximize the movie's probability of profitability. More formally, the problem is to, for $\forall v_i \in C_M$, construct

candidate cast member set C'_{M_i} , where $\forall v'_i \in C'_{M_i}$ is a candidate that could replace $v_i \in C_M$. In other words, we want to construct a set of candidate casts $C'_M = \text{Combine}\left(C'_{M_1}, C'_{M_2}, \dots, C'_{M_{|C_M|}}\right)$, where $\text{Combine}(\bullet)$ is a function that returns the non-overlapping combinations of elements from each respective candidate actor set. This process yields a total of $|C'_M| = \prod_{i=1}^{|C_M|} |C'_{M_i}|$ candidate cast sets. As can be seen, the number of possibilities is combinatorially large. For example, for a movie with five actors, if we can identify ten candidates for each of the five actors, then the total number of possible cast sets to be evaluated is 10^5 .

To reduce the search space, we required each of the candidates in the new cast to be similar to each of the original cast members in three ways: star power, gender, and age. Star power (*SP*) serves as a proxy for the amount of pay needed to hire an actor (measured by cumulative total gross). Gender and age help to ensure that the role is aptly represented (e.g., the role of a teenage girl in a movie is usually played by a young female actress). Age was approximated by tenure, the length of time a particular actor has been appearing in movies. We also required that i' has appeared in a movie in the past y years to ensure, as best as possible, that candidates i' are still actively pursuing work in the movie industry. Specifically for our experiments, the following similarities need to exist between i and i' : (1) $SP_{i'} \in [SP_i * 0.90, SP_i * 1.05]$, so that the cost of hiring a replacement i' is similar to i ; (2) $\text{gender}_{i'} = \text{gender}_i$; (3) $\text{tenure}_{i'} \in [\text{tenure}_i - 7, \text{tenure}_i + 7]$, so that the i' 's age is close to that of i ; and (4) $\text{activity}_{i'} \in [M_{\text{release}} - 3, M_{\text{release}}]$ to ensure that i' has been active during the past three years.

In our experiments, for each candidate actor set $s \in C'_M$, we computed cast-dependent features (i.e., changeable features) and appended them to non-cast-dependent features (i.e., unchangeable features that are the same for each $s \in C'_M$). This created a unique movie M_s , to which we can apply our predictive model that returns the probability of M_s having $ROI \geq 24$ percent. After predicting each M_s using our random forest classifier, we can select $M_s^* = \max\{P(M_s), \forall s \in C'_M\}$, where $P(\cdot)$ denotes the probability of being profitable from the predictive model, and M_s^* is the movie, starred in by candidate cast set s , that is most likely to earn $ROI \geq 24$ percent. The process of computing all cast-dependent features for all $|C'_M|$ candidate cast sets is time-consuming. A more comprehensive exploration of this process requires a more efficient procedure for finding M_s^* .

We conducted experiments on two different movies—*Abandon*, a 2002 PG-13-rated thriller, and *Captain Corelli's Mandolin*, a 2001 R-rated war-romance movie. Both were originally unprofitable, according to our $ROI \geq 24$ percent definition, and their inability to be profitable was correctly predicted by our models. Also, both movies had three first-billed actors. Table 12 summarizes the experiment results for the two movies including the original cast, the number of candidate cast members for each original cast member, and the “ideal cast” that would maximize the probability of being profitable. Had the two movies selected the “ideal cast,” the predicted probability of having $ROI \geq 24$ percent would have increased from 0.07 to 0.35 for *Abandon*, and 0.08 to 0.47 for *Captain Corelli's Mandolin*.

Table 12. Cast Selection Results for Two Movies

| Movie | Cast Info. | Actor 1 | Actor 2 | Actor 3 |
|-----------------------------------|-----------------|------------------|-----------------|-----------------|
| <i>Abandon</i> | Original Cast | Zooey Deschanel | Katie Holmes | Benjamin Bratt |
| | # of Candidates | 2 | 4 | 10 |
| | 'Ideal Cast' | Carrie-Anne Moss | Kate Bosworth | David Schwimmer |
| <i>Captain Corelli's Mandolin</i> | Original Cast | Nicolas Cage | Christian Bale | Penelope Cruz |
| | # of Candidates | 7 | 7 | 10 |
| | 'Ideal Cast' | Antonio Banderas | Michael Nyqvist | Christina Ricci |

Limitations

Admittedly, our study has limitations. Similar to past studies, the profit we calculated is based on estimated production budget and reported box office revenue. However, the true profit of a movie may be obscured by certain accounting practices. Teasing out the effect of such practices is very challenging without sensitive accounting data and such an endeavor is beyond the scope of this study. In addition, for many movies, box office revenue is only one of the sources of income. For example, Disney's animation movies often gain a significant part of their revenues from the sale of movie merchandise, such as clothing and toys. Some movies may also rely heavily on the sale and rental of DVDs. However, capturing these non-box-office revenues is more difficult because they can continue accumulating many years after the release of a movie.

Conclusions and Future Work

In this study, we proposed a movie investor assurance system (MIAS) to aid movie investment decisions at the early stage of movie productions. MIAS learns from freely available historical data derived from various sources, and tries to predict movie success based on profitability. The system's performance improves greatly with novel features that we proposed for the first time, including "who" is in the cast, "what" a movie is about, "when" a movie will be released, and "hybrid" features that match these features. In addition to predicting whether a movie is worth investing in, our research also informs prescriptive analytics—MIAS not only allows "what-if" analysis in order to experiment with what increases the chance of profitability, but also can be the basis of cast recommendation to select profit-maximizing cast members for a movie.

Besides movie investors, our system can also be helpful for other stakeholders in the movie industry who care about the possible financial success of a movie, such as

cinemas that want to decide whether to air a movie. Moreover, the framework of MIAS, as well as the features we extracted for MIAS, can also be applied to other creative work, which often requires a team of contributors, whose content can be described with texts, and for which timing is important, such as research papers, grant proposals, operas, and so on.

The research highlights the power of data analytics in building information systems that support business decision making. The outcomes potentially have theoretical implications as well. Our regression analysis revealed the effects of key factors of movie profitability. Some findings contradict the results of previous studies. For example, the importance of directors for movie profitability was highlighted. Our new methods of quantifying factors suggested by past theoretical studies (e.g., actor star power, team expertise, and team diversity) also worked better in the context of predicting profitability. We hope these findings will inspire future theoretical research in areas such as marketing, creative work, and team performance.

There are also several directions for future research. For example, as we have matched “what” with “when” and “what” with “who,” it would be interesting to match “who” with “when” to capture whether the popularity of an actor or director is on the rise or declining. Another interesting future direction for research would be to collect full-length scripts of a large number of movies and to analyze entire scripts, instead of the plot synopses. Movie scripts can provide more details on movies’ content, as well as novel features, such as script cadence. We also intend to add more features to our model, including those that more definitively address consumer spending power, such as external economic indices, as well as those that take into account the types of movies that are most suited to certain times of the year (i.e., is it best to release Christmas-themed movies at Christmas time?). Analyzing how successful, or well-known, the source of an adapted movie is could also contribute to the prediction of movie profitability. Our method of prescribing cast members can be improved by adding more realistic cast-selection criteria and reducing the computational complexity.

REFERENCES

1. Abbasi A.; Zahedi F. M.; Zeng D.; Chen Y.; Chen H.; and Nunamaker, J.F. Jr. Enhancing predictive analytics for anti-phishing by exploiting website genre information. *Journal of Management Information Systems*, 31, 4 (2015), 109–157.
2. Aggarwal, C.C.; Chen, C.; and Han, J. The inverse classification problem. *Journal of Computer Science and Technology*, 25, 3 (2010), 458–468.
3. Apala, K.R.; Jose, M.; Motnam, S.; Chan, C.C.; Liszka, K. J.; and de Gregorio, F. Prediction of movies box office performance using social media. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining*. Niagara Falls: IEEE Computer Society, 2013, pp. 1209–1214.
4. Asur, S., and Huberman, B.A. Predicting the future with social media. In *Proceedings of the 2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology*. Toronto: IEEE Computer Society, 2010, pp. 492–499.

5. Baimbridge, M. Movie admissions and rental income: The case of James Bond. *Applied Economics Letters*, 4, 1 (1997), 57–61.
6. Blei, D.M.; Ng, A.Y.; and Jordan, M.I. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3, 1 (2003), 993–1022.
7. Boccardelli, P.; Brunetta, F.; and Vicentini, F. What is critical to success in the movie industry? A study on key success factors in the Italian motion picture industry. *Dynamics of Institutions and Markets in Europe*, 46, 4 (2008), 1–22.
8. Bozdogan, Y. The determinants of box office revenue: a case based study: Thirty, low budget, highest ROI films vs. thirty, big budget, highest grossing Hollywood films. MATHesis, University of Paris, 2013.
9. Burt, R.S. Structural holes and good ideas. *American Journal of Sociology*, 110, 2 (2004), 349–399.
10. Burt, R.S. *Structural Holes: The Social Structure of Competition*. Cambridge, MA: Harvard University Press, 1992.
11. Chi, C.L.; Street, W.N.; Robinson, J.G.; and Crawford, M.A. Individualized patient-centered lifestyle recommendations: An expert system for communicating patient specific cardiovascular risk information and prioritizing lifestyle options. *Journal of Biomedical Informatics*, 45, 6 (2012), 1164–1174.
12. Craney, T.A., and Surles, J.G. Model-dependent variance inflation factor cutoff values. *Quality Engineering*, 14, 3 (2002), 391–404.
13. Cui, G.; Wong, M.L.; and Wan, X. Cost-sensitive learning via priority sampling to improve the return on marketing and CRM investment. *Journal of Management Information Systems*, 29, 1 (2012), 341–374.
14. Elberse, A. The power of stars: Do star actors drive the success of movies? *AMA Journal of Marketing*, 71, 4 (2007), 102–120.
15. Eliashberg, J.; Hui, S.; and Zhang, Z. Assessing box office performance using movie scripts: A kernel-based approach. *IEEE Transactions on Knowledge and Data Engineering*, 26, 11 (2014), pp. 2639–2648.
16. Eliashberg, J.; Hui, S.K.; and Zhang, Z.J. From story line to box office: A new approach for green-lighting movie scripts. *Management Science*, 53, 6 (2007), 881–893.
17. Eliashberg, J.; Jonker, J.J.; Sawhney, M.S.; and Berend, W. MOVIEMOD: An implementable decision-support system for prerelease market evaluation of motion pictures. *Marketing Science*, 19, 3 (2000), 226–243.
18. Gopinath, S.; Chintagunta, P. K.; and Venkataraman, S. Blogs, advertising and local market movie box office performance. *Management Science*, 59, 12 (2013), 2635–2654.
19. Guimera, R.; Uzzi, B.; Spiro, J.; and Amaral, L.A.N. Team assembly mechanisms determine collaboration network structure and team performance. *Science*, 308, 5722 (April 2005), 697–702.
20. Hevner, A.R.; March, S.T.; Park, J.; and Ram, S. Design science in information systems research. *MIS Quarterly*, 28, 1 (2004), 75–105.
21. Kuhn M., and Johnson, K. *Applied Predictive Modeling*. New York: Springer, 2013.
22. Lash, M.T.; Fu, S.; Wang, S.; and Zhao, K. Early prediction of movie success: What, who, and when. In N. Agarwal, K. Xu, and N. Osgood (eds.), *Proceedings of the 2015 International Conference on Social Computing, Behavioral-Cultural Modeling, and Prediction*. Washington, DC: Springer, 2015, pp. 345–349.
23. Lash, M.T.; Lin, Q.; Street, W.N.; and Robinson, J.G. A budget-constrained inverse classification framework for smooth classifiers. arXiv preprint, 2016. <https://arxiv.org/abs/1605.09068>
24. Lutter, M. Creative Success and Network Embeddedness: Explaining Critical Recognition of Film Directors in Hollywood, 1900–2010 (July 9, 2014). MPIfG Discussion Paper 14/11. Available at SSRN: <https://ssrn.com/abstract=2464150>
25. Magni M.; Angst C.M.; and Agarwal R. Everybody needs somebody: The influence of team network structure on information technology use. *Journal of Management Information Systems*, 29, 3 (2012), 9–42.
26. Meiseberg, B.; and Ehrmann, T. Diversity in teams and the success of cultural products. *Journal of Cultural Economics*, 37, 1 (2013), 61–86.

27. Meiseberg, B.; Ehrmann, T.; and Dormann, J. We don't need another hero: Implications from network structure and resource commitment for movie performance. *Schmalenbach Business Review*, 60, 1 (2008), 74–99.
28. Mestyan, M.; Yasseri, T.; and Kertész, J. Early prediction of movie box office success based on Wikipedia activity big data. *PLoS ONE*, 8, 8 (January 2013). <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0071226>
29. Motion Picture Association of America (MPAA). 2015 Theatrical market statistics. 2015. http://www.mpa.org/wp-content/uploads/2016/04/MPAA-Theatrical-Market-Statistics-2015_Final.pdf
30. Parimi, R.; and Caragea, D. Pre-release box-office success prediction for motion pictures. In *Proceedings of the Ninth International Conference on Machine Learning and Data Mining in Pattern Recognition*. New York: Springer Berlin Heidelberg, 571–585.
31. Prat, N.; Comyn-Wattiau, I.; and Akoka, J. A taxonomy of evaluation methods for information systems artifacts. *Journal of Management Information Systems*, 32, 3 (2015), 229–267.
32. Sharda, R.; and Delen, D. Predicting box-office success of motion pictures with neural networks. *Expert Systems with Applications*, 30, 2 (2006), 243–254.
33. Sinha, A.P., and May, J.H.; Evaluating and tuning predictive data mining models using receiver operating characteristic curves. *Journal of Management Information Systems*, 21, 3 (2004), 249–280.
34. Taylor, P.; Simonoff, J.S.; and Sparrow, R. Predicting movie grosses: Winners and losers, blockbusters and sleepers. *CHANCE*, 13, 2 (2014), 15–24.
35. Uzzi, B., and Spiro, J. Collaboration and creativity: The small world problem. *American Journal of Sociology*, 111, 2 (2005), 447–504.
36. Vany, A.D.E., and Walls, W.D. Uncertainty in the movie industry: Does star power reduce the terror of the box office? *Journal of Cultural Economics*, 23, 4 (1999), 285–318.
37. Wallace, W.T.; Seigerman, A.; and Holbrook, M.B. The role of actors and actresses in the success of films: How much is a movie star worth? *Journal of Cultural Economics*, 17, 1 (1993), 1–27.
38. Walls, W.D. Modeling movie success when nobody knows anything: Conditional stable distribution analysis of film returns. *Journal of Cultural Economics*, 29, 3 (2005), 177–190.
39. Zaheer, A., and Soda, G. Network evolution: Structural holes. *Administrative Science Quarterly*, 54, 1 (2007), 1–31.
40. Zhang, W., and Skiena, S. Improving Movie Gross Prediction through News Analysis. In *Proceedings of the 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology*. Milan: IEEE Computer Society, 2009, pp. 301–304.
41. Zhao, K.; Wang, X.; Yu, M.; and Gao, B. User recommendation in reciprocal and bipartite social networks: A case study of online dating. *IEEE Intelligent Systems*, 29, 2 (2014), 27–35.
42. Zhao, K.; Yen, J.; Ngamassi, L.M.; Maitland, C.; and Tapia, A.H. Simulating inter-organizational collaboration network: a multi-relational and event-based approach. *Simulation*, 88, 5 (2011), 617–633.

Copyright of Journal of Management Information Systems is the property of Taylor & Francis Ltd and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.