

FORECASTING AUDIENCE OF MOTION PICTURES CONSIDERING COMPETITIVE ENVIRONMENT

¹SEONYEONG PARK, ²TAEГУ KIM

¹MS. Student, Hanbat National University,
Department of Industrial and Management Engineering, Daejeon, Korea

²Corresponding author, Assistant Professor, Hanbat National University,
Department of Industrial and Management Engineering, Daejeon, Korea

E-mail: ¹paksy725@gmail.com, ²taegu.kim@hanbat.ac.kr

ABSTRACT

Film industry is one of the most prospective cultural business sectors and attracts much of public attention. However, investment on individual film titles is notorious for its high risk level, which has raised the necessity of research on the influential factors for box office revenue. Despite of abundant extant studies, those attempts have been mostly limited to a certain kind of factors and the competitive environment is even hardly explored. In this regard, we investigate the significance and relative contribution of a wide variety of possible influential factors including distribution power and competition. In this study, we devised various new variables to reflect competition environment and categorized independent variables into several groups in order to compare their influences. Results showed that all various information reflected in four variable groups contribute to box office revenue in both estimation and forecasting. Between mathematical models, SVM models exhibit better result than linear models as expected. Regarding competition variables, the number of competitors is found to be more effective compared to their distribution power.

Keywords: *Forecasting, Movie, Box office, Machine Learning, Competition, Linear Regression*

1. INTRODUCTION

Today, the film industry is one of the steadily growing industries worldwide [1], [2]. Korea is also experiencing rapid growth in the movie industry, and the annual frequency of visiting theaters is ranked at the highest level among all countries [3]. Due to this importance as an industry and the inherent high risk as a cultural product, various studies have been conducted to maximize the return on investment [1], [2], [4]–[6].

A variety of studies has extensively explored the motion picture industry over the past 30 years. Most studies have focused on analyzing the influential factors and predicting the box office revenues [1], [2], [4]–[10]. Variables commonly investigated in those studies are the basic characteristics of movie titles such as the number of screenings, distribution parameters such as screen number, distributor, production, genre, film ratings, nationality, seasonality, director, star, and sequel [1], [2], [5], [7], [8], [10]–[15]. In addition to these factors,

reviews from critics and audiences were also covered [1], [7]–[11], [16].

The most prevalent analytic method chosen in previous studies was linear regression analysis, which selected significant variables among the various factors and analyzed how they affected the dependent variables. On the other hand, movie titles were categorized according their characteristics to figure out the difference in influential factors. Further, optimal managerial decision for release strategy, screening frequency and period were suggested based on forecasting results with the chosen factors [1], [2], [5], [7]–[9], [11]–[15], [17].

In addition to the linear regression analysis, diffusion models and machine learning techniques were also adopted to improve forecasting accuracy and interpreting diffusion patterns [2], [6], [17]–[23]. Among various machine learning algorithms such as artificial neural network, decision tree, and k-nearest neighbors, support vector machine (SVM, hereafter) is the most remarkable approach in this area for its excellence in fitness and forecasting performances [2], [6], [17]–[23].

Though many studies have analyzed influential factors on audiences, there have been few studies dealing with various factors comprehensively. Therefore, in this study, we intend to forecast the gross number of audiences by introducing competition variables in addition to the other variables covered in previous studies[16], [24]–[28]. Moreover, variables are categorized into several groups and their influences are compared.

Main contributions of this study can be viewed in two ways as follows. First, forecasting accuracy is improved by including many variables found significant in previous studies. Particularly, in order to establish a model considering the competitive environment, which has not been explored enough in previous research, various competition variables are devised as numbers and distribution power of competitors. Moreover, categorized variable groups are evaluated by comparing their contributions to estimation fitness. Second, forecasting models are diversified and improved by employing both linear regression and SVM. Two methods contribute through their strengths respectively. Linear regression is utilized to select significant variables and SVM is adopted for its superiority in nonlinear problems.

This paper is composed as follows. Section 2 describes the data acquisition and preprocessing process and explains the variables. Section 3 describes the estimation and forecasting results. Finally, Section 4 discusses conclusions and suggests future research.

2. DATA

In this study, daily box office data of Korean movie market for 700 days from August 1, 2014 to June 30, 2016 were used. The collected data includes data on rank, sales, market share, number of audiences, number of screens, number of screenings, nationality, director, actor, distributor, and production.

Originally, there were two attributes related to distribution power: the number of screens and the number of screenings (or screening frequency). In addition, the product of these two values was also considered as candidate for distribution power. To avoid multi-collinearity problem, three separate single linear regressions were conducted. Consequently, screening frequency was selected with the highest adjusted R^2 value.

From the 478 titles ever ranked in daily top 10 lists, 23 titles released before August 2014 were excluded and total of 455 films were finally selected for analysis. Next, preprocessing was conducted: adding new variables expected to contribute to the forecasting such as director power and competitive environment, log-transformation of dependent and several independent variables related to number of audiences, screens and screening frequencies.

Table 1 shows the variables preprocessed to be analyzed. Seasonality variable is allocated as one if the target movie was released in holiday season (January to February or July to August) or as zero otherwise. Director power (*dir_power*) is defined as the number of audiences of the last title of the same director prior to the target movie.

Table 1 : Description for variables

Variable	Group	Log
Total Audience	Dependent	Log
Frequency on Opening Day	Distribution Power	Log
Distributor		-
Nationality		-
Ratings	Movie	-
Genre	Characteristic	-
Seasonality		-
Director Power		Log
# of titles with same nationality		-
# of titles with same genre	Competition	-
# of titles with same ratings	-Numbers	-
# of titles with same distributor		-
# of screens with same nationality		Log
# of screens with same genre	Competition	Log
# of screens with same ratings	-Distribution	Log
# of screens with same distributor		Log

Meanwhile, competition variables are devised to reflect the competitive environment around the target title and measured with numbers and distribution powers of competitors. Competitors were defined as concurrent titles in daily top 10 box office list similar to the target movie. Similarity is evaluated with four criteria: nationality, genre, ratings, and distribution.

Number of competitors is counted as of the previous day of the opening day and distribution power of competitors is determined as the sum of the screens allocated to competitors on the opening day.

Also, to compare the influence of the variables, we classified them into four groups according to their characteristics. First, as revealed in many previous studies, distribution power that are considered to be the most influential factor separately forms a group. Next, variables related to the identity of the target movie, such as genre or ratings, are included in the movie characteristic variable group. Finally, the competition variables, which are the main concern of this study, were divided into two groups according to the method of measuring the intensity of competition as mentioned above.

On the other hand, test dataset was separately prepared to evaluate forecasting accuracy of the model. Titles in the test dataset are ‘Magnificent 7’, ‘Gosanja’, and ‘Cafe Society.’

3. EXPERIMENT AND RESULT

3.1 Experiment Design

The purpose of this study is to analyze the influence of various variables including distribution power, movie characteristics and competitive environment on the number of audiences. To this end, we compare the estimation and prediction performance of each model with five combinations of variable groups as shown in Figure 1 below.

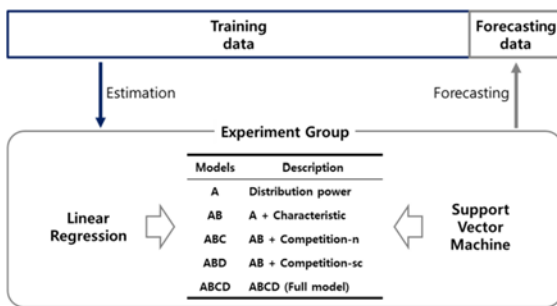


Figure 1. Experiment design

Overall, the experiment is done by comparing the performance of the five models. Each model consists of different variable groups. First, model A with only distribution power variable becomes the baseline. Next, in order to compare the performance of each variable group in the manner of incremental contribution, other models are constructed by adding each variable group to the baseline model. Consequently, characteristic variables were incorporated in model AB and thereafter, whereas competition variables are added afterwards. With

this structure, the degree of contribution of each variable group to audience demand can be grasped through the difference in performance between models. For example, the performance difference between model A and AB can be interpreted as the effect of variable group B.

Furthermore, to improve the accuracy of the model, support vector machine as a nonlinear machine learning algorithm is employed as well as simple linear regression. In the case of regression analysis, only significant variables were selected for each model through stepwise variable selection method.

3.2 Estimation Result

The overall performance shown in Figure 2 and Table 2 indicates that the SVM performs better than the linear regression, and fitness measures are improved with the addition of variables in general. In particular, the prediction accuracy is significantly improved by the variable group B in the linear regression model. For the competitive environment, the numbers of competitors are more influential than their distribution powers.

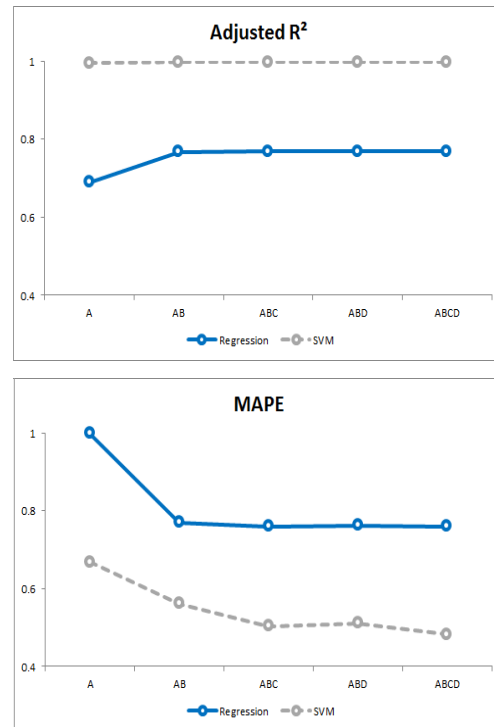


Figure 2. Estimation result

Table 3: Summary of estimation results

Group Index	Group/Subgroup	Variable Description	A	AB_selected	ABC_selected	ABD_selected	ABCD_selected	
		(Intercept)	4.367***	3.807***	3.957***	3.854***	3.957***	
A	Distribution power	Frequency on Opening Day	1.177***	1.208***	1.202***	1.213***	1.202***	
		CGV		0.453**	0.429**	0.453**	0.429**	
		CJ		0.305*	0.366*	0.292**	0.366*	
		CINE GURU						
		The Walt Disney		0.959***	1.080***	0.933***	1.080***	
		Isu C&E						
		20th Fox Korea		0.730***	0.822***	0.735***	0.822***	
		Little Big Pictures		-0.435**	-0.389**	-0.416**	-0.389**	
		Distributor	Lotte Ent.					
			Megabox					
			NEW			0.287		0.287
			Others					
			Pan Cinema					
			Pop Ent.		-0.529**	-0.495**	-0.548*	-0.495**
	Showbox		0.762***	0.770***	0.749***	0.770***		
	UPI		0.224	0.286**	0.234	0.286**		
	Wanna Bros.		0.758***	0.870***	0.740***	0.870***		
B	Nationality	Others						
		USA		-0.336***	-0.311***	-0.332***	-0.311***	
	Ratings	ALL		0.566***	0.444***	0.622***	0.444***	
		12+		0.168		0.304*		
		15+				0.182		
	Genre	Animation		0.485***	0.513***	0.467***	0.513***	
		Comedy		-0.215	-0.296**	-0.229	-0.296**	
		Drama		0.287***	0.474***	0.280***	0.474***	
		Horror			-0.309		-0.309	
		Melodrama						
Others								
	Thriller			-0.323		-0.323		
Season	Seasonality		0.160*	0.158*	0.178***	0.158*		
Director power	Director power							
C	Competition -Numbers	Nationality_n						
		Genre_n			-0.055*	-0.055*		
		Rating_n						
		Distributor_n						
D	Competition - Distribution	Nationality_sc						
		Genre_sc						
		Rating_sc			-0.033*			
		Distributor_sc						

Asterisks indicate the levels of significance ***: 0.01, **: 0.05, *: 0.1

Table 2: Estimation results

Model	Regression		SVM	
	Adj.R ²	MAPE	Adj.R ²	MAPE
A	0.6891	99.85%	0.9946	66.80%
AB	0.7675	77.01%	0.9970	56.16%
ABC	0.7690	76.02%	0.9973	50.37%
ABD	0.7687	76.29%	0.9972	51.11%
ABCD	0.7690	76.02%	0.9975	48.22%

Figure 3 is the scattered plot of estimation results with linear regression and SVM. Estimated points from SVM are much closer to the actual data compared to the linear model results.

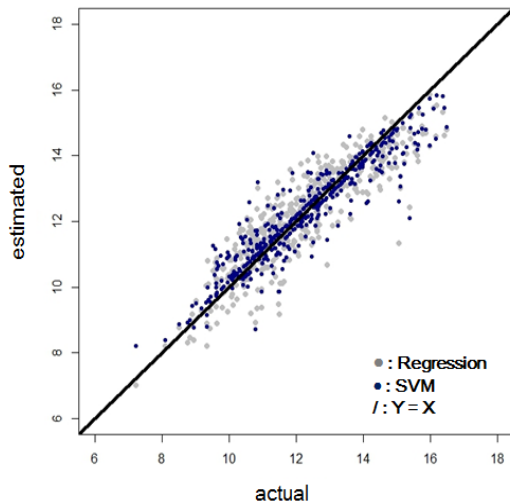


Figure 3. Scattered plot of estimation result

Table 3 summarizes variable selection results for linear regression models. For five different models (A, AB, ABC, ABD, ABCD), coefficient estimates are written in each row of the table. Director power variable was expected to be significant, yet is found to be not influential. By contrast, many movie characteristic variables such as distributor, rating, nationality, seasonality, and genre are tested as significant. Results with seasonality and ratings hint that a bigger potential audience group leads to a higher box office score.

Above all, the distribution power is found to be significant for all models, which corresponds to previous studies [11], [21], [23], [29]. It is plausible that the many screens on the opening day means the high expectation level and provides more opportunities for audiences.

This is also related to the distributors, because bigger the distribution company is, more screening frequency the movie has. Actually, large distributors such as Disney, 20th Fox Korea, Wanna Bros., and Showbox are significant across all models, while medium distributors such as CINE GURU and Pan Cinema and small distributors classified as "others" are not significant in any case. Even some small distributors (Pop Ent., Little Big Pictures) turned out to have negative impacts on the box office scores.

Investigation on director power concluded that the previous success of a director would not guarantee the success of the following work. This insignificance can be also found in the previous studies [12], [30]

Results on other movie characteristic factors, variables in group B, genre and seasonality do not raise any new arguments against existing facts. Influential power of genre is not decisive as only several dummy variables are significant [31], [32]. On the contrary, seasonality is consistently significant across all models corresponding to existing arguments [13], [33]. However, it is notable that titles imported from US fall behind Korean films.

We tested a variety of competition variables, yet only one cases for each variable group C and D were determined as influential. Among the number of competitors, the number of movies with the same genre was found to inhibit the box office. On the other hand, the distribution power of movies with the same ratings reduces the target movie's box office score. Between these two factors, it seems that the number of competitors with same genre is more meaningful than the other because it is the only competition variable selected in the full model (ABCD). This result can be interpreted as people choose movies on a title basis and does not know nor consider how many screens they occupy.

In the manner of incremental contribution for the fitness, movie characteristic variables are much better than the competition variables. For both linear regression and SVM, adding group B improved the fitness far more than adding group C or D, while the overall trend of incremental improvement are clearly demonstrated in the MAPE graph with SVM.

3.3 Forecasting Result

Overall, forecasting performance was better than estimation. Based on the full model (ABCD),

forecasting results with linear regression and SVM were 54% and 17% which are much more favorable than estimation errors of 76% and 42% respectively. Table 4-6 show the detail results.

Table 4: Forecasting Results (Magnificent 7)

Model	Regression		SVM	
	Forecasting	APE	Forecasting	APE
Actual	921,317	-	921,317	-
A	725,939	21.2%	678,568	26.3%
AB	471,960	48.8%	715,180	22.4%
ABC	481,122	47.8%	976,861	6.0%
ABD	481,942	47.3%	793,033	13.9%
ABCD	481,122	47.8%	904,195	1.9%

Table 5: Forecasting Results (Gosanja)

Model	Regression		SVM	
	Forecasting	APE	Forecasting	APE
Actual	974,225	-	974,225	-
A	1,030,824	5.8%	1,469,996	50.9%
AB	1,807,364	85.5%	689,120	29.3%
ABC	1,870,194	92.0%	672,939	30.9%
ABD	2,048,095	110.2%	620,791	36.3%
ABCD	1,870,194	92.0%	647,977	33.5%

Table 6: Forecasting Results (Café Society)

Model	Regression		SVM	
	Forecasting	APE	Forecasting	APE
Actual	126,385	-	126,385	-
A	112,786	10.8%	61,557	51.3%
AB	87,789	30.5%	155,413	23.0%
ABC	97,898	22.5%	116,482	7.8%
ABD	88,843	29.7%	128,993	2.1%
ABCD	97,898	22.5%	107,240	15.1%

In the case of ‘Magnificent 7’, the full model (ABCD) produced a very close to the actual gross audience. However, adding variable group B over the baseline model rather hinders the forecasting performance. Forecasting accuracies were most disappointing for ‘Gosanja.’ This is mainly because the movie captured much attention from audience at its release, yet failed to meet the expectation level and consequently made a relatively poor score. Linear regression models present the most favorable forecasting results for “Café society,” where the contribution of each variable group is similar in estimation and forecasting.

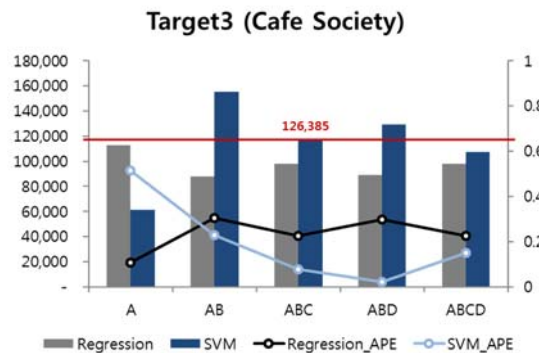
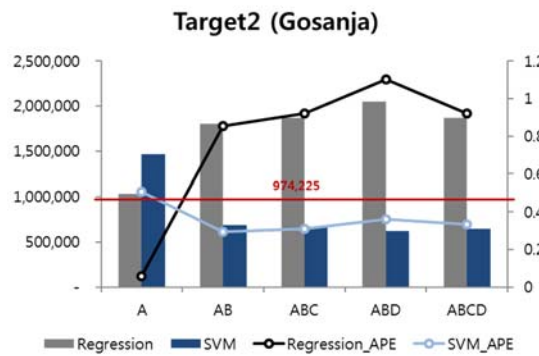
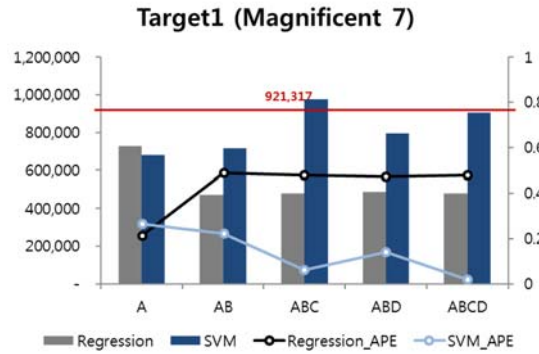


Figure 4. Forecasting results

As Figure 4 depicts above, the overall pattern of result does not coincide with our conjecture of effects of variable groups on forecasting performance. Actually incorporating characteristic information (group B) rather hampers forecasting accuracies in most of cases, which means that the forecasting targets are much different from the titles in estimation dataset.

Further, the performance difference between SVM and linear models is more distinct in forecasting results compared to estimation results.

4. CONCLUSION

4.1 Summary and Discussion

The importance of the movie market and the high risk of the movie titles have been the source of various studies on the influential factors. In this study, we attempted to investigate the influence of each variable group by adding variables to reflect competition environment and grouping variables, as well as including various factors found in previous studies. In addition, to improve the predictive power, a nonlinear machine learning technique was used as well as linear regression method.

Four variable groups all helped to improve the estimation and prediction performance and SVM showed better performance than linear regression model. Though the characteristic factors such as genre and ratings contributed the most, it did rarely aid in the forecasting accuracy. Meanwhile, for the competitive factors, the number of competitors was found to be more crucial than their distribution power.

Regression analysis showed that the variables chosen to be significant differed slightly from those of the public recognition in general. Awareness level of the director is considered as a relevant factor yet was found not to be significant in this study.

4.2 Limitation and Future Research

In this study, there is a limit of the inconsistency between the estimation and the prediction results, which is mainly due to the insufficient number of titles for forecasting. Therefore, it is necessary to expand the test dataset and revise the results.

Next, in the case of informal data, results may vary according to the method of definition, so research to reflect informal information such as director power more accurately is required. Furthermore, analysis on the differences in influential factors across clusters according to nationality, genre, ratings and size should be carried out.

ACKNOWLEDGMENTS

This research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education, Science and Technology (NRF-2017R1C1B1010854)

REFERENCES:

- [1] Kim, S.Y., Im, S.H. and Jung, Y.S., "A Comparison Study of the Determinants of Performance of Motion Pictures: Art film vs. Commercial film," *The Korea Contents Society*, Vol. 10, No. 2, 2010, pp. 381-393.
- [2] Jeong, H.Y. and Yang, H.J., "Predicting Financial Success of a Movie Using Multiple Regression Analysis," *Proceedings of the Korean Society of Computer Information Conference*, Vol. 21, No. 2, 2013, pp.275-278.
- [3] Korean Film Council, "Korea Film Industry in 2016," 2016.
- [4] Ji, Y.G., "Individuals' Demand for Movies and Theaters' Choice of Movies on the Screen," Ph.D. dissertation, Department of Economics & Finance, Hanyang University, Seoul, Korea, 2015.
- [5] Wang, B.S. and Chon, B.S., "Determinants of Movie Success from Foreign Countries in Korea," *Journal of the Korea Contents Association*, Vol. 16, No. 2, 2016, pp. 96-105.
- [6] Park, Y.E., Kim, S.H., Park, H.J. and Rhee, D.K., "Exploratory Study on the Factors Influencing the Profitability of Korean Movies," *Korean Management Review*, Vol. 39, No. 2, 2010, pp. 459-488.
- [7] Kim, B.S., "Comparison of Factors Predicting Theatrical Movie Success: Focusing on the Classification by the Release Type and the Length of Run," *Korean Journal of Journalism & Communication Studies*, Vol. 53, No. 1, 2009, pp. 257-287.
- [8] Park, S.H. and Jung, W.K., "The Determinants of Motion Picture Box Office Performance: Evidence from Movies Released in Korea, 2006-2008," *Journal of Communication Science*, Vol. 9, No. 4, 2009, pp. 243-276.
- [9] Lee, Y.H., Chang, B.H. and Park, K.W., "An Exploratory Study for Comparing Factors Affecting Box Office Performances between Countries: Focusing on Performances of U.S. Movies in South Korea and U.S.," *Journal of Communication Science*, Vol. 7, No. 1, 2007, pp. 185-222.

- [10] Kim, D.S., "Integrated Research on Motion Picture Success Factor: Based on Information Available Period," M.S. thesis, Division of Management Engineering, Korea Advanced Institute of Science and Technology, Daejeon, Korea, 2006.
- [11] Duan, W., Gub, B. and Whinston, A.B., "The dynamics of online word-of-mouth and product sales an empirical investigation of the movie industry," *Journal of Retailing*, Vol. 84, No. 2, 2008, pp. 233-242.
- [12] Elberse, A. and Eliashberg, J., "Demand and supply dynamics for sequentially released products in international markets: the case of motion pictures," *Marketing Science*, Vol. 22, No. 3, 2003, pp. 329-354.
- [13] Litman, B.R., "Predicting success of theatrical movies: an empirical study," *Journal of Popular Culture*, Vol. 16, No. 4, 1983, pp. 159-175.
- [14] Lovallo, D., Clarke, C. and Camerer, C., "Robust analogizing and the outside view: two empirical tests of case-based decision making," *Strategic Management Journal*, Vol. 33, No. 5, 2012, pp. 496-512.
- [15] Ravid, S.A., "Information, blockbusters, and stars: a study of the film industry," *The Journal of Business*, Vol. 72, No. 4, 1999, pp. 463-492.
- [16] Hababou, M., Amrouche, N. and Jedidi, K., "Measuring Economic Efficiency in the Motion Picture Industry: a Data Envelopment Analysis Approach," *Cust. Needs Solut.*, vol. 3, no. 3-4, pp. 144-158, 2016.
- [17] Abel, F., Diaz-Aviles, E., Henze, N., Krause, D., and Siehdnel, P., "Analyzing the blogosphere for predicting the success of music and movie products," In *Proceedings of the 2010 international conference on advances in social networks analysis and mining*, Odense, Denmark, 2010, pp. 276-280.
- [18] Ghiassi, M., Lio, D. and Moon, B., "Pre-production forecasting of movie revenues with a dynamic artificial neural network," *Expert Systems with Applications*, Vol. 42, No. 6, 2015, pp. 3176-3193.
- [19] Delen, D., Sharda, R. and Kumar, P., "Movie forecast guru: a web-based DSS for Hollywood managers," *Decision Support Systems*, Vol. 43, No. 4, 2007, pp. 1151-1170.
- [20] Zhang, L., Luo, J. and Yang, S., "Forecasting box office revenue of movies with BP neural network," *Expert Systems with Applications*, Vol. 36, No. 3, 2009, pp. 6580-6587.
- [21] Kim, T.G. and Hong, J.S., "Identifying the Diffusion Patterns of Movies by Opening Strength and Profitability," *Journal of the Korean Institute of Industrial Engineers*, Vol. 39, No. 5, 2013, pp. 412-421.
- [22] Kang, J.H., Park, C.H., Do, H.R. and Kim, S.B., "Methodologies in Forecasting Box office Revenue using Datamining Techniques," *Proceedings on 2014 Spring Joint Conference of KIIE and KORMS*, Busan, Korea, 2014, pp. 142-154.
- [23] Kim, T.G., Hong, J.S. and Kang, P.S., "Box office forecasting using machine learning algorithms based on SNS data," *International Journal of Forecasting*, Vol. 31, No. 2, 2015, pp. 364-390.
- [24] Yang, J.H., Kim, W. J., Amblee, N. and Jeong, J. S., "The heterogeneous effect of WOM on product sales: why the effect of WOM valence is mixed?," *Eur. J. Mark.*, vol. 46, no. 11-12, pp. 1523-1538, 2012.
- [25] Delre, S.A., Panico, C. and Wierenga, B., "Competitive strategies in the motion picture industry: An ABM to study investment decisions," *Int. J. Res. Mark.*, vol. 34, no. 1, pp. 69-99, 2013.
- [26] Yu, Y. and Chen, H. "Measuring Social Media Success: The Case of Facebook Marketing in the Motion Picture Industry," in *Pacific Asia Conference on Information Systems (PACIS) Proceedings*, 2015, pp. 1-9.
- [27] Divakaran, P. K. P., Palmer, A., Søndergaard, H. A. and Matkovskyy, R., "Pre-launch Prediction of Market Performance for Short Lifecycle Products Using Online Community Data," *J. Interact. Mark.*, vol. 38, pp. 12-28, 2017.
- [28] Prieto-Rodriguez, J., Gutierrez-Navratil, F. and Ateca-Amestoy, V., "Theatre allocation as a distributors strategic variable over movie runs," *J. Cult. Econ.*, vol. 39, no. 1, pp. 65-83, 2015.

- [29] Kim, T., Hong, J., and Koo, H., “Forecasting Box-Office Revenue by Considering Social Network Services in the Korean Market,” *J. Teknol.*, vol. 64, no. 2, pp. 97–101, 2013.
- [30] Wen, K.H., and Yang, C.Q., “Determinants of the Box Office Performance of Motion Picture in China-Indication for Chinese Motion Picture Market by Adapting Determinants of the Box Office (Part I),” *J. Sci. Innov.*, p. 25, 2011.
- [31] Sawhney, M.S., and Eliashberg, J., “A parsimonious model for forecasting gross box-office revenues of motion pictures,” *Mark. Sci.*, vol. 15, no. 2, pp. 113–131, 1996.
- [32] Neelamegham, R., and Chintagunta, P., “A Bayesian Model to Forecast New Product Performance in Domestic and International Markets,” *Mark. Sci.*, vol. 18, no. 2, pp. 115–136, 1999.
- [33] Zhang, W., and Skiena, S., “Improving Movie Gross Prediction through News Analysis,” *Proceedings of the 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology - Volume 01*. IEEE Computer Society, pp. 301–304, 2009.