

Visual Analytics for the Prediction of Movie Rating and Box Office Performance

Mennatallah el Assady
Christian Rohrdantz

Daniel Hafner
Fabian Fischer

Michael Hund
Svenja Simon

Alexander Jäger
Tobias Schreck

Wolfgang Jentner
Daniel A. Keim *

University of Konstanz, Germany

1 INTRODUCTION

This paper describes our solution to the IEEE VAST 2013 Mini Challenge 1¹. The task of the challenge was to create a visual and interactive tool to predict the popularity of new movies in terms of viewer ratings and ticket sales for the opening weekend in the U.S. The data usage was restricted by the challenge organizers to data from the Internet Movie Database (IMDb)² and a predefined set of Twitter³ microblog messages. To tackle the challenge we designed a system together with an analysis workflow, combining machine learning and visualization paradigms in order to obtain accurate predictions. In Section 2 we describe the machine learning components used within the analysis workflow. Next, in Section 3, we describe where and how the human analyst is enabled to enhance the prediction with her/his world knowledge. Finally, Section 4 concludes the paper providing an evaluation of the prediction accuracy with and without human intervention.

2 MACHINE LEARNING

In order to predict the performance of ratings and box office takings for upcoming movies, it makes sense to rely on data from past movies. For example, if movies from certain directors or with certain actors have been successful in the past, it is to be expected that also their future movies will be successful. It has been shown that machine learning models, like neural networks, can support the prediction in this way [1]. We experimented with different models for predicting the movie viewer rating, which we trained and tested based on IMDb data from past movies. As input we took into account movies with same cast and crew members or genres. In a 10-fold cross validation test, neural network predictors performed best. For the prediction of box office takings, we applied multinomial regressions of different orders. The input parameters to the regressions were the movie budget and runtime, as those were the only reliable numerical variables available across all movies. Another line of research has shown that social media messages, such as Tweets, can potentially also be exploited in a beneficial way when predicting movie performance [4]. Yet, when predicting the performance of upcoming movies as requested by the VAST 2013 Mini Challenge 1 certain limiting and biasing factors have to be taken into account. First, the challenge demands the prediction of the rating a movie achieves on the opening weekend only. In contrast to that, the IMDb data contains ratings that have established over a longer period of time. Thus, the characteristics of the training data do not match those of the data that is to be predicted. Another related issue is that box office takings from the distant past are probably less meaningful for the prediction due to economic factors such as infla-

tion and developments in ticket pricing and purchasing power. Second, the challenge organizers restricted the sample of Tweets that were allowed to be used to a predefined set. We found out that this set was hardly representative and sometimes also contained many Tweets not related to the corresponding movie. Past research suggests that there are correlations between numerical characteristics that could be derived from the Tweets, such as number of Tweets or Tweet sentiments [4], and the rating performance. However, such correlations did not appear within our restricted Tweet sample and the set of movies to be predicted. Thus, the Tweets were not useful for generating an automatic prediction. Third, the automatic methods lack the integration of world knowledge. Especially in the prediction of box office takings external factors not contained in the data may have a strong impact. For example, the number of cinemas in which a movie is shown, the coincidence of holidays, the weather on the opening Weekend, whether a movie is shown 2D or 3D or both, and which other movies are released on the same weekend running in competition for spectators. For the prediction of viewer ratings important external factors are, for example, whether the movie is based on a book, whether it is a sequel, and the publicity for the movie spread in news or through web channels.

In order to account for these biases it is required to integrate the human into the analysis loop. Still, an automatic prediction provides an indication for what a realistic value or range of values for the final prediction could be. In the next section we will detail on how the human analyst can enhance the automatic prediction.

3 INTERACTIVE ADJUSTMENT

As mentioned, mere automatic methods fall short of incorporating all possibly available factors influencing the prediction. The human analyst has to be integrated into the analysis workflow in order to contribute world knowledge and interpretations of social media content. Our solution is twofold: First, the analyst interactively decides on what s/he considers to be the most useful input for the machine learning. We name this pre-learning interaction phase. Second, after the machine learning process has finished the analyst is enabled to tune the results. We name this post-learning interaction phase. For both phases we offer different interactive visual displays. For the sake of brevity, only some fundamental steps will be described in the following paragraphs. Further details are given in our challenge submission⁴ and video⁵, which are both available online.

Pre-learning interaction phase First, we provide a graph-based visualization that reveals details on the social media content relating to a certain movie. The graph-structure shows co-occurrences of different persons, concepts, and attributes and also reflects sentiments. The graph visualization is created using VISONE⁶. The structure of the graph is generated as follows: *Nodes* represent different types of keywords (names of actors, adjectives, verbs, nouns, and #hashtags). Each of these types is mapped to

*firstname.lastname@uni-konstanz.de

¹<http://boxofficevast.org/>

²<http://www.imdb.com/>

³<https://twitter.com/>

⁴<http://bib.dbvis.de/uploadedFiles/MooVisSummaryFinal.pdf>

⁵youtu.be/XhJDPa9FNck

⁶<http://visone.info/html/about.html>

