

Received May 22, 2016, accepted June 12, 2016, date of publication June 14, 2016, date of current version July 7, 2016.

Digital Object Identifier 10.1109/ACCESS.2016.2580911

# A Peek Into the Future: Predicting the Popularity of Online Videos

**SHUXIN OUYANG<sup>1</sup>, CHENYU LI<sup>2</sup>, (Student Member, IEEE), AND XUEMING LI<sup>2</sup>**

<sup>1</sup>School of Information and Communication Engineering, Beijing University of Posts and Telecommunications, Beijing 100876, China

<sup>2</sup>Beijing University of Posts and Telecommunications, Beijing 100876, China

Corresponding author: S. Ouyang (ouyangshuxin@gmail.com)

**ABSTRACT** Predicting the popularity of online videos is an important task for the service design, advertisement placement, network management, and so on. In this paper, we tackle the challenge head-on by casting the popularity prediction problem into two consecutive tasks: online video future popularity level prediction and online video future view count prediction. We first predict the future popularity levels of online videos, based on a rich set of features and effective classification technique. Then, according to the popularity level transitions, we build specialized regression models to predict the future view count values. We validate our approach on the exhaustive dataset of a leading online video service provider in China, namely, Youku. The experimental results show that comparing with two state-of-the-art baseline models, our proposed method can significantly decrease the relative prediction errors of 32.25% and 19.82%, respectively. At last, we also discuss the model setup and feature importance of our method. We believe our work can provide direct help in practical for the interested parties of online video service, such as service providers, online advisers, and network operators.

**INDEX TERMS** Online video service, video popularity prediction.

## I. INTRODUCTION

Among the various kinds of Web 2.0 services, online video service is currently the dominating one on the Internet. As a white paper of Cisco System notes [1], in terms of bytes the video traffic accounted for 70% of all the Internet traffic in 2015, and will be up to 82% by 2020. Every second, nearly a million minutes of video content crosses the network, and it will take an individual over 5 million years to watch all the videos across the Internet each month. With such large volume of video content competing the limited user attention and time, it is not surprised that ecosystem of online videos turns into the asymmetric *rich-get-richer* effect. A small fraction of the videos attracts most of the user interest, whereas the vast majority are barely noticed by users [2].

Hence, predicting the future popularity of online video content, which is measured by view count, is of great importance for a number of contexts. For the service providers, such knowledge can help them design more effective information services, such as video ranking, video searching and recommendation schemes [3]. For the online advertisers, by identifying the next rising star of online videos, they can better place their advertisements and estimate the revenue in advance [4]. And for the network operators, with the

estimations of incoming video views, delivery network resources (e.g. bandwidth and cache servers) can be proactively allocated to avoid the potential bottlenecks [5]. Moreover, in a quite different context, the future popularity of online videos is of great interest for the opportunistic communications among mobile devices [6]. In such resource-constrained environments (e.g. limited bandwidth and storage), predicting hot videos can provide help for the strategies of content delivering, caching and replicating.

In this paper, we study the video popularity of Youku ([www.youku.com](http://www.youku.com)), a leading online video service provider in China. Our work is based on the data of 200,714 videos collected from Youku for 30 consecutive days. With these data, We tackle the challenge of predicting future popularity of online videos by proposing a two-stage prediction method. First, we estimate the future popularity level of a video based on a rich set of features and effective classification technique. Then, according to the transition from the observed popularity level to the predicted popularity level, we build specialized regression models to predict the precise value of future view count for the video. We test our method on a real-world dataset and compare the prediction performance with two state-of-the-art baseline models of online video

popularity prediction. The experimental results show that our method leads to significant reductions in the prediction errors, reaching up to 32.25% and 19.82% over the baseline models, respectively. We further analyze the potentials and limitations of different detectors in the popularity level prediction. In addition, we investigate the importance and utility of each feature and feature group, and shed light on the key impact factors related to the video popularity. The main contributions of our work are summarized as follows:

- We release a large-scale, long-term and up-to-date online video analysis database to the public.<sup>1</sup>
- We provide in-depth analysis of the importance of different features in the online video future popularity level prediction.
- We propose a two-stage method to predict the future popularity of online video.
- The proposed two-stage online video future popularity prediction model gets state-of-the-art results and outperforms conventional methods.

The rest of this paper is organized as follows: We firstly outline the related work in Section II. Then in Section III, we briefly describe our dataset and provide a characterization of the video popularity distribution. The details of popularity prediction, including the model, the experiment, and the discussion, are provided in Section IV, Section V and Section VI. At last, we conclude our paper in Section VII.

## II. RELATED WORK

### A. POPULARITY LEVEL PREDICTION OF ONLINE CONTENT

Jamali and Rangwala [7] extracted several user-based and comment-based features for Digg stories, and trained different classification methods to predict the future popularity level with those features collected during the first ten hours of a Digg story. Lee *et al.* [8] introduced a widely used survival analysis method (Cox proportional-hazards regression) into the problem of online content popularity prediction, and detected whether a DPreview or MySpace thread would still be popular (with comments more than a threshold) after a certain period of time. Tsagkias *et al.* [9] addressed the problem of predicting the comment volume for online news story as two consecutive classification tasks. They first identified whether a story would receive a comment, and then based on the results they predicted the comment volume to be low or high. Vallet *et al.* [10] collected cross-system features from YouTube and Twitter websites, and used gradient boosted decision tree method to classify the YouTube videos into popular, viral, or both categories. Roy *et al.* [11] also extracted information from Twitter, and associated popular topics to YouTube videos. They proposed an algorithm to describe the social transfer and further identified the YouTube videos that would experience sudden bursts of popularity.

### B. POPULARITY NUMBER PREDICTION OF ONLINE CONTENT

Szabo and Huberman [2] observed a strong linear correlations between the log-transformed long-term popularity and log-transformed early popularity. Based on this characteristic, they proposed a simple log-linear model to predict the future popularity of online content. This method was verified on various kinds of online content: YouTube videos [2], Digg stories [2], online news articles [12]–[14] and etc. In [2], the daily view count increments were treated equally. It is not difficult to notice that treating daily increments equally is not accurate. Pinto *et al.* [15] further modified the log-linear model, and proposed a multivariate regression model by assigning different weights to the daily increments in view count during the early observation period. After that, Vasconcelos *et al.* [16] used an ordinary least squares (OLS) multivariate linear regression model based on user, venue and content features, to predict the micro-review popularity of a location based social networking platform, Foursquare. Li *et al.* [17] analyzed how the popularity of Youku videos evolved over time, and used multivariate regression models with different parameters to predict the future popularity of the videos with different popularity evolution patterns. In addition to these regression based methods, there are also efforts towards building models to predict online content popularity using other techniques such as reservoir computing [18], time series analysis [19] and hidden Markov model (HMM) [20].

Our study complements these existing works by considering the popularity level prediction and the future view count prediction conjointly. We cast the prediction of video popularity into a classification task and a regression task, and propose a two-stage prediction method. Moreover, we introduce the notion of popularity level transition into building specialized regression models, and eventually gain significant improvements of the prediction performance.

## III. PRELIMINARIES

### A. DATASET

The data used in this paper were collected from a leading online video service provider in China, namely Youku. Youku archives more than 500 million monthly active users and 800 million daily video views [21]. It provides a comprehensive type of online video service, including both user-generated content (UGC) and video-on-demand (VoD). To study the popularity of Youku videos, we crawled and tracked the **video meta-data** and the **user meta-data** for a set of videos via Youku Open API [22]. The **video meta-data** includes the static and dynamic properties of a video, such as: title, description, category, duration, published time, tags, stream type, copyright type, public type, source name, view count, comment count, favorite count, up count, down count, user ID and etc. And the **user meta-data** includes the statistic information of the video uploader, such as: registration time, account status, video count, video view count, favorite count,

<sup>1</sup>[https://github.com/lichenyu/Datasets/tree/master/Youku\\_Popularity\\_151212\\_151221](https://github.com/lichenyu/Datasets/tree/master/Youku_Popularity_151212_151221)



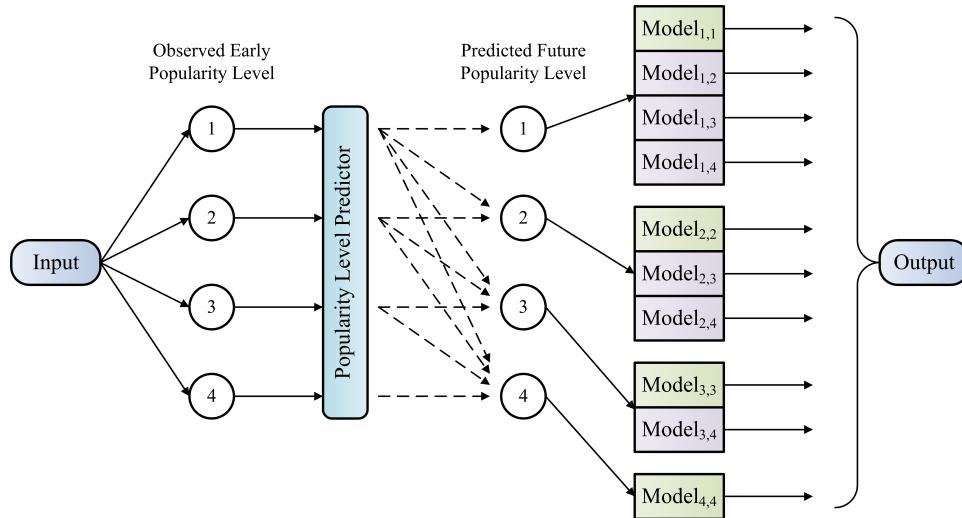


FIGURE 2. Overview of the prediction system.

**User Statistic Features** capture the account information (e.g. registration time, whether a VIP), publication activities (e.g. the number of uploaded videos) and social influence (e.g. the number of followers) of the video uploader.

**Content Topic Features:** The content topic of a video is closely related to its popularity [24]. In our analysis, we use the presences of certain representative terms in video tags and titles to serve as content topic features. More specifically, for each popularity level, we extract the top 500 terms that most frequently appear in the tags and titles (in the training set). Then, the common terms that appear across popularity levels, as well as the words without actual meanings such as “of” and “very”, are removed from the word sets. At last, we keep the top 100 words for each popularity level as the word list, and use bag-of-word (i.e. frequency) features for the videos in the test set as the content topic features.

**Textual Analysis Features** capture the syntactic and sentiment information of the textual content of the video. Texts in Youku are mainly written in Chinese. We utilize *jieba* [25] to extract the part-of-speech (PoS) features of the video title. We also use *snownlp* [26] to compute the sentiment score of the video title, description, and tags. A sentiment score ranges from 0 to 1, indicating the text from very negative to very positive sentiment.

**Historical Popularity Features** measure the popularity evolution of a video during the early observation period, including the daily increment of view count, the ratio of daily increment over total view count in the observation period, and other video popularity metrics (e.g. comment, rating and favorite) the video archives at the end of the observation period. Moreover, we describe the popularity evolution pattern of the video by utilizing the frequency of popularity burst in the early observation period. In our analysis, if the daily increment ratio ( $\frac{\text{daily increase in view count}}{\text{total view count}}$ ) is larger than three times of the average ratio (i.e. 1/7),

we consider the video experiences a popularity burst on that day.

For the prediction of popularity level, with those 69 features, we exploit the random forest [27] technique as the classifier.

## B. FUTURE VIEW COUNT PREDICTION

In future view count prediction, there are two kind of situations: 1) popularity level stays the same between reference date and observed date, 2) popularity level changes between reference date and observed date.

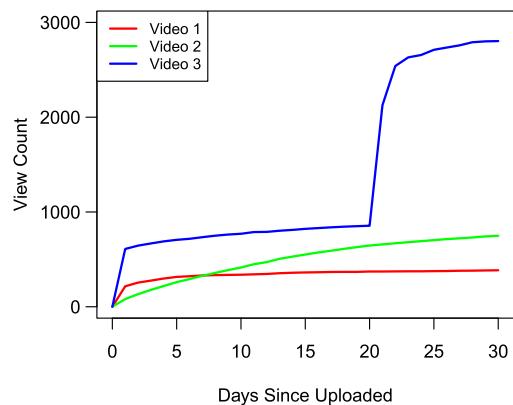


FIGURE 3. Popularity evolutions over time for three different videos.

### 1) SITUATION 1

If a video’s predicted popularity level on the reference date remains the same as its observed popularity level on the indicator date, the most important indicator for future popularity should be the popularity evolution pattern during the early observation period (described by the increases in view count). Take the popularity evolutions of three videos in our dataset for examples, as shown in Fig. 3. For video 1







