

CROWDSOURCED TIME-SYNC VIDEO TAGGING USING SEMANTIC ASSOCIATION GRAPH

Wenmian Yang¹, Na Ruan¹, Wenyuan Gao¹, Kun Wang², Wensheng Ran³, Weijia Jia^{1*}

¹Shanghai Jiao Tong University, China,

²Nanjing University of Posts and Telecommunications, China

³Tokyo Institute of Technology, Japan

{sdq11111, gaowenyuan}@sjtu.edu.cn, {naruan, jia-wj}@cs.sjtu.edu.cn,
kwang@njupt.edu.cn, ranwensheng@ecei.tohoku.ac.jp

ABSTRACT

Time-sync comments reveal a new way of extracting the on-line video tags. However, such time-sync comments have lots of noises due to users' diverse comments, introducing great challenges for accurate and fast video tag extractions. In this paper, we propose an unsupervised video tag extraction algorithm named Semantic Weight-Inverse Document Frequency (SW-IDF). SW-IDF first generates corresponding semantic association graph (SAG) using semantic similarities and timestamps of the time-sync comments. Then it clusters the comments into sub-graphs of different topics and assigns weight to each comment based on SAG. This can clearly differentiate the meaningful comments with the noises. In this way, the noises can be identified, and effectively eliminated. Extensive experiments have shown that SW-IDF can achieve 0.3045 precision and 0.6530 recall in high-density comments; 0.3800 precision and 0.4460 recall in low-density comments. It is the best performance among the existing unsupervised algorithms.

Index Terms— video tagging, crowdsourced time-sync comments, semantic association graph, keywords extraction

1. INTRODUCTION

Recently, many people spend their time on on-line video sites for news and entertainment. The wide applications of on-line (or live) videos raise great challenges in fast and accurate videos searching techniques. There are many efforts made on automatic video tagging techniques based on text retrieval [1, 2], but they can only provide video-level tags [3]. In other words, the description of the video tags are not specific enough and timestamps are not matched with tags. Fortunately, a new type of comments, i.e., time-sync video comments appear on video websites like Youku (www.youku.com), AcFun (www.acfun.tv) and Bilibili (www.bilibili.com) in China, and NicoNico (www.nicovideo.jp) in Japan.

In this paper, we focus on extracting video tags from time-sync comments to obtain more accurate and specific tags. Time-sync video comment is an interactive comment form, where users can make their comments synchronized to a video's playback time. These comments are saved and will be displayed on the video screen for other users in the future.

Also, they are crowdsourced data [4, 5], which convey information involving the content of the current video, feelings of users or replies to other comments. Therefore, time-sync video comments mining can be served as a new way for video tag extraction.

As far as we know, only Wu *et al.* [3] do some work on extracting video tags from time-sync comments, which inspires our work. They use statistics and topic model to build Temporal and Personalized Topic Modeling (TPTM). However, their approach needs the user id of each comment, which is difficult to obtain accurately because of the privacy-protecting policy. Furthermore, their approach does not consider the semantic of comments, leading some of the video content-independent noises can not be processed.

To tackle these problems, we propose a graph-based algorithm named Semantic Weight-Inverse Document Frequency (SW-IDF) to generate time-sync video tags automatically. The time-sync video comments have some features distinguished from the common short text, so the main idea in our algorithm is based on these features. Specifically, time-sync video comments have the following features: (1) *Semantic relevance*. Abundant video semantic information is contained that describes both local and global video contents by selecting the time interval of the timestamp. (2) *Real-time*. Time-sync video comment is synchronous to the real-time content of the videos. Users may produce different topics when it comes to the same video contents. (3) *Interdependence*. Latter comments usually depend on the former ones, which means the latter comments have a semantic association with the foregoing. (4) *Noise*. Plenty of noises and internet slang are involved in comments, which makes trouble for tag extraction. Therefore, how to distinguish the weight of each comment and consequently identify high-impact comments and noises by utilizing the features above is a major challenge.

According to our observation of plenty of time-sync video comments, the noises have few semantic relationships with other time-sync comments while some high impact comments have mass semantic relationships with others in a period of time. Based on this, we can reduce the impact of noises by clustering the semantic similar and time-related comments, and identify high impact comments by their semantic relationships. Specifically, we first generate corresponding semantic association graph (SAG) using semantic similarities

and timestamps of the time-sync comments. Then we treat the time-sync comments as vertices in the graph and cluster them into different topics by community detection theory [6]. The weight of each comment is assigned based on the degrees of comment in. This can clearly differentiate the meaningful comments with noises. Moreover, we gain the weight of each word by combining semantic weight (SW) and inverse document frequency (IDF) which is similar to the TF-IDF algorithm and then video tags are extracted in an automatic way.

The main contributions of our paper are as follows:

- 1) We propose a novel graph-based Semantic Weight-Inverse Document Frequency (SW-IDF) algorithm, which can extract both local and global video tags in an unsupervised way by mining time-sync comments.
- 2) We build Semantic Association Graph (SAG) to cluster the comments into sub-graphs of different topics. The method takes the features of time-sync comments into account, and effectively reduce the impact of noises.
- 3) We evaluate our proposed algorithm with real-world datasets on mainstream video-sharing websites and compare it with classical keyword extraction methods. The results show that SW-IDF outperforms baselines in both precision and recall of video tag extraction.

2. RELATED WORK

In this section, we talk about the related work from three aspects.

Time-sync video comment has an increasing number of researches emphasize text mining on it. In addition to Wu *et al.*'s work [3], there are also some other applications based on time-sync comments. Xian *et al.* [7] extract highlight of video clips by mining the time-sync comments in a simplified model. Lv *et al.* [8] propose a T-DSSM to represent comments into semantic vectors, and video highlights are recognized by semantic vectors in a supervised way. Wu and Ito [9] investigate the correlation between emotional comments and popularity of videos and He *et al.* [10] propose a model to predict the popularity of videos. Although their applications are different from ours, their idea of mining the time-sync comments inspired our work.

Tag/keyword extraction is another work related to ours. At present, three unsupervised keyword extraction methods are available. The first one is based on word frequency statistics, such as TF-IDF. The second method depends on the relationship between the sentences, such as textrank [11], which is a graph-based ranking model. And the last one is according to topic model. It brings document-topic and topic-word distribution together by simulating document generation process. Blei *et al.* [12] propose the Latent Dirichlet Allocation(LDA) model, the most representative model. To better deal with short text situation, Yan *et al.* propose the BTM [13], which models the generation of word co-occurrence patterns (i.e. biterns) in the whole corpus directly. Yin and Wang propose the Gibbs Sampling algorithm for the Dirichlet Multinomial Mixture model [14, 15] for short text clustering and keyword extraction. However, methods above cannot deal with the texts that are full of noises.

Semantic similarity of time-sync comments is a critical issue in our paper. There are mainly two kinds of approaches to measure the similarity of documents. One is based the similarity of the words in sentences. The representations of this approach are propose by Li *et al.* [16] on unsupervised learning and Socher *et al.* [17, 18] on supervised learning. Considering that time-sync comments contain a mass of newborn internet slangs, it is difficult to obtain accurate results in this way. The other one is based on the sentence vector. Texts are converted to vector as first, and distance of the vectors is calculated as the similarity between sentences. The topic model such as LDA, and embedding model such as word2vec [19] are the representations of this method. Since embedding model offers much denser feature representation, embedding based similarity computation is better in this paper.

3. OUR ALGORITHM

In this section, we will first introduce the construction of Semantic Association Graph (SAG) for time-sync comments with their semantic similarity. Then we cluster the comments into subgraphs of different topics, generating the weight of each comment by an out-in degree iterative algorithm based on SAG. Finally, keywords are extracted as video tags from time-sync comments automatically.

3.1. Preliminaries and Graph Construction

In this subsection, we will construct the semantic association graph and provide the definition of the attributes in the graph.

Since comments appear in chronological order, they can only affect the upcoming comments rather than prior comments. We use a directed graph to describe the relationships between comments and construct the semantic association graph (SAG). In SAG, the vertices are comments and the edges reflect their semantic association in a topic. Let G denote the directed graph. It is represented by $G = (V, E)$, where V and E are the sets of nodes and edges. Specifically, $V = \{v_1, v_2, \dots, v_N\}$, $E = \{e_1, e_2, \dots, e_M\}$, where N is the number of nodes in V , and M is the number of edges in E . In our algorithm, each comment i has a timestamp t_i , denoting the post time of the time-sync comment. $t_{v_1} < t_{v_2} < \dots < t_{v_N}$. Each comment in our algorithm has one exact topic. For vertex v_i , $v_i.S$ is used to describe the set that contain the vertices which have the same topic as v_i and $|S|$ is used to express the number of vertices in set S . We use the domain to describe the attributes of edges. For each edge e_k , $e_k.x$ and $e_k.y$ are two vertices that are linked by edge e_k where $t_{e_k.x} < t_{e_k.y}$. The weight of edge i is described as $e_i.w$. Besides, $e_{u,v}$ also describes the edge with vertices u and v where $t_u < t_v$. Next we will provide the definition of edge weights.

As we mentioned in Section 2, an embedding based method word2vec is selected to calculate the semantic similarity between each pair of comments because there are abundant training data of time-sync comments. To simplify the algorithm, we calculate the mean vector of each word in the sentence as the sentence vector. The dimensionality of each vector is 300. Therefore, the semantic similarity between comment a and b is calculated by the cosine angle between

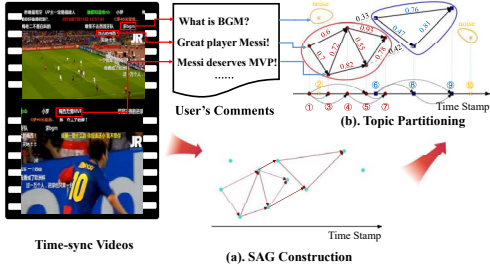


Fig. 1. An example of SAG Construction

vectors:

$$Sim(a, b) = \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|}$$

Besides, the greater the timestamp interval between two comments, the less likely they are in the same topic. So we use the exponential function to express the decay of comment associations:

$$Delay(a, b) = exp^{-\gamma_t \cdot (t_b - t_a)}$$

where γ_t is an attenuation coefficient.

Combining the semantic similarity and the time decay, the weight of edge i that link vertices u and v is defined by i that link vertices u and v is defined by

$$e_{i.w} = \begin{cases} Sim(u, v) \cdot Delay(u, v) & t_u < t_v \\ 0 & t_u > t_v \end{cases}$$

Empirically derived threshold, comments which edge weight are less than 0.3 are intuitively too semantic dissimilar. So when $e_{u,v}.w < 0.3$, we can set $e_{u,v}.w = 0$ and delete this edge.

For a more intuitive description, an example of SAG construction is shown in Fig. 1 (a). In a UEFA Champions League video, user A made the comment 1 as “Great player Messi!” when he saw the goal. Then user B responded with “Messi deserves MVP!” as the comment 3. User C makes a comment “What is the BGM ?” as comment 2 to ask the background music, which deviates the video content. So the comment made by user C has no semantic association with other comments, while comments proposed by user A and B have a semantic edge.

3.2. Topic Partitioning

In this subsection, we will partition the topic of each comment according to the semantic relationships in SAG. In our algorithm, the time-sync comment that has the similar semantics and similar timestamps should belong to the same topic. So the mean weight of edges in intra-topic is large while the mean weight of edges that link different topics is small, which satisfies community detection theory.

In the beginning, each comment belongs to a unique set that only contains itself to achieve the objective. That is,

for each v_i , $v_i.S = \{i\}$. Then edges in set E are sorted by descending order of weight. The new edge set $E' = \{e'_1, e'_2, \dots, e'_k, \dots, e'_M\}$ is obtained, where $e'_1.w > e'_2.w > \dots > e'_M.w$. We process each edge from e'_1 to e'_M . For each edge e'_k , we use S_1 and S_2 to represent the sets $e'_k.x.S$ and $e'_k.y.S$. We merge S_1 and S_2 when $S_1 \neq S_2$ and

$$\frac{\sum_{e_p.x, e_p.y \in S_1 \cup S_2} e_p.w}{(|S_1| + |S_2|) \cdot (|S_1| + |S_2| - 1)/2} > \rho$$

where ρ is the threshold of intra-cluster density. In this paper, disjoint-set (union-find set) algorithm [20] is used to merger the sets efficiently. When all the edges are solved, comments with high semantic similarity are merged into a topic, and the intra-cluster density of each subgraph is higher than the threshold. An example about topic partitioning is shown in Fig. 1 (b). The SAG constructed in Fig. 1 (a) is partitioned into two topics and several noises. The comment “Great player Messi!” and “Messi deserves MVP!” belong to the same topic, while the comment “What is the BGM ?” is identified as a noise.

3.3. Weight Distribution and Tag Extraction

Algorithm 1 EXTRACTING TAGS BY SW-IDF

Input Semantic Association Graph

Output Tags of video

- 1: sort E by descending order of $e_i.w$, obtain E'
 - 2: **for** $i = 1$ to M **do**
 - 3: set $e'_i.x.S$ as S_1 , $e'_i.y.S$ as S_2
 - 4: **if** $\frac{\sum_{e_p.x, e_p.y \in S_1 \cup S_2} e_p.w}{(|S_1| + |S_2|) \cdot (|S_1| + |S_2| - 1)/2} > \rho$ **then**
 - 5: merge S_1 and S_2
 - 6: **end if**
 - 7: **end for**
 - 8: Calculate the influence matrix $M_{N \times N}$ using Eq.(2)
 - 9: **for** $i = 1$ to N **do**
 - 10: $I_i^0 = 1$
 - 11: Calculate the popularity of comment i using Eq.(1)
 - 12: **end for**
 - 13: **for** $k = 1$ to T **do**
 - 14: **for** $i = N$ downto 1 **do**
 - 15: Calculate I_i^{2k-1} using Eq.(3)
 - 16: **end for**
 - 17: **for** $i = 1$ to N **do**
 - 18: Calculate I_i^{2k} using Eq.(4)
 - 19: **end for**
 - 20: **end for**
 - 21: Calculate the SW-IDF of each word using Eq.(5)
 - 22: Select words with max SW-IDF as video tags
-

We partition the topic in section 3.2 and get the topic of each comments. In this section, we will attribute weight to each comment according to the influence of its topic and the

relationship in the semantic graph.

The weight of a comment is affected by its topic popularity, so we define the popularity of the comment i as:

$$P_i = \frac{|v_i \cdot S|}{\sqrt[K]{|S_1| \cdot |S_2| \dots |S_K|}} \quad (1)$$

where S_i is the i -th topic in semantic association graph, and K is the total number of topics in semantic association graph. Obviously those topics with fewer comments are more likely to be noises and have less weight. According to Eq.(3), noises will have small values of popularity.

Within the topic, a comment which affects more comments and is affected by less comments should have higher weight. In order to quantitatively measure the weight of the comment in a topic, we design a graph iterative algorithm. An influence matrix $\mathbb{M}_{N \times N}$ is established at first to semantic relations within each topic. For the elements in the matrix,

$$m_{i,j} = \begin{cases} e_{i,j} \cdot w & v_i \cdot S = v_j \cdot S \\ 0 & v_i \cdot S \neq v_j \cdot S \end{cases} \quad (2)$$

We use $I_{i,k}$ to denote the influence value of i -th comment after k iterations. For each comment i , $I_{i,0} = 1$ initially. Then in the k -th turn of iteration, there are two steps as follows:

$$I_{i,2k-1} = I_{i,2k-2} + \sum_{j=i+1}^n m_{i,j} \cdot I_{j,2k-1} \quad (3)$$

and

$$I_{i,2k} = \frac{I_{i,2k-1}}{I_{i,2k-1} + \sum_{j=1}^{i-1} m_{j,i} \cdot I_{j,2k}} \quad (4)$$

In the $(2k-1)$ -th iteration, we increase the influence value of comment i based on the influential value of its affecting comments. We know that a comment will only affect the comments lagging behind it, so the comments can be processed from v_N down to v_1 . That is, before we process comment i , all the comments j that $t_j > t_i$ have been processed. In the $(2k)$ -th iteration, we reduce the influence value of comment i based on the influence value of the comments which affect comment i . Contrary to the $(2k-1)$ -th iteration, we process the comments from v_1 to v_N in the $(2k)$ -th iteration. The influential value of the comments in Fig. 1 is shown in Fig. 2. After 20 iterations, all comments converge to the interval $[0, 1]$. In topic 1, the influence value of comment 3 and comment 4 is bigger than that of comment 6 and comment 7, meeting our expectation.

To combine the popularity and the influence value, the comment weight i is obtained by

$$W_i = P_i \cdot I_i^T$$

where T is the number of turns of iterations and depends on the number of nonzero elements in matrix $\mathbb{M}_{N \times N}$. Therefore, the weight of each word is formulated as below:

$$SW - IDF_i = \sum_j W_j \cdot IDF_i \quad (5)$$

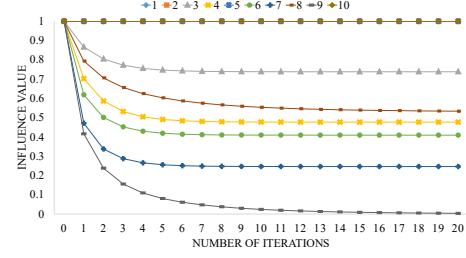


Fig. 2. The influence value of comments in Fig. 1

where j denotes the comment that word i appears and IDF_i is the inverse document frequency as defined in TF-IDF method. We extract words with the highest SW-IDF value as video tags. After the above steps, those words which appear in the comments that are popular and have high impact will be extracted as tags. The complete algorithm is shown in Algorithm 1.

4. EXPERIMENTAL STUDY

In this section, we verify the effectiveness of our proposed method by comparing with four unsupervised methods of keyword extraction. The datasets are crawled from AcFun. We provide necessary parameters for our model at first and then analyze the performance of our method on video tag extraction.

4.1. Experimental Setup and Datasets

We crawl time-sync comments from a Chinese time-sync comments video website AcFun. To be specific, totally 227,780 comments are collected randomly from music, sports, and movie, 126,146 comments for the training set and 101,634 comments for the test set. The test set is divided into two parts: high-density comments and low-density comments by density (the number of comments per second). We artificially make video tags of 120 videos as the standard. More details included Length (second), Number of comments, Density (comments per second) and the number of videos about test set are shown in Table 1.

Table 1. Data Description Table

	Length	Comments number	Density	Video number
High density	7475	41,556	5.5593	89
Low density	79008	60,078	0.7604	31

In our algorithm, two parameters need to be determined, *i.e.*, the threshold of intra-cluster density ρ , and the attenuation coefficient γ_t . The ρ controls the accuracy of topic clustering. The γ_t is the attenuation coefficient of the interval between time-sync comments, which controls the value of the edge weights in the graph.

To determine ρ , we randomly choose some comments from both low-density and high-density comments in the test set. Then we artificially partition the comments that should

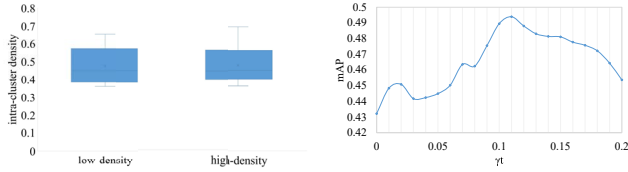


Fig. 3. Box diagram of intra-cluster density **Fig. 4.** The effect of attenuation coefficient γ_t

belong to the same topic and calculate the intra-cluster density of each topic. The box diagram of intra-cluster density is shown in Fig. 3. The minimum value of intra-cluster density in both low-density and high-density is nearly 0.36, so ρ is initiated to 0.36 in our experiment.

To determine γ_t , we randomly choose 30 tagged videos from the test set, adjusting γ_t and observe the mAP (mean average precision) of video tags generated by our model. Results are shown in the Fig. 4. γ_t gains better performance in the range of 0.1 to 0.15 and gets optimal performance at 0.11. In fact, when $\gamma_t = 0$, the semantic association graph is independent of time; when $a = +\infty$, all weights of edge equal to 0, and our model is equivalent to TF-IDF. Besides, the number of iterations T also needs to be determined. We find in the experiment that the graph with hundreds of vertices is convergent after nearly 20 interactions and that with thousand vertices is convergent after approximately 50 interactions.

4.2. Results

In this section, we use the test set that are described in Section 4.1 to compare our algorithm with existing methods.

Regrettably, we cannot choose Wu et al.’s [3] as baseline, because both video sites AcFun and Bilibili have already hidden the ID of users when posting comments to protect privacy, which is necessary in Wu et al.’s method. To evaluate the performance of the proposed video tag extraction algorithm, we compare our method with 4 unsupervised keyword extraction methods, *i.e.*,

- (1) A classical keyword extraction algorithm TF-IDF.
- (2) A graph-based text ranking model, textrank [11], which is inspired by PageRank. The method is denoted as “TX”.
- (3) A topic model based algorithm, Biterm Topic Model [13], which is the improvement of LDA [12] for short texts. The method is denoted as “BTM”.
- (4) A collapsed Gibbs Sampling algorithm for the Dirichlet Process Multinomial Mixture model [14, 15], which has good performance when dealing with short texts. This method is denoted as “GSDPMM”

Table 2. Comparison of different methods on video tag extraction of the top 10 candidate tags with high density comments.

Method	Prec	Recall	F-score	mAP
TF-IDF	0.2674	0.5735	0.3648	0.4224
TX	0.2427	0.5205	0.3310	0.3696
BTM	0.2337	0.5012	0.3188	0.3094
GSDPMM	0.2445	0.5094	0.3302	0.3374
SW-IDF	0.3045	0.6530	0.4153	0.4853

Table 3. Comparison of different methods on video tag extraction of the top 10 candidate tags with low density comments.

Method	Prec	Recall	F-score	mAP
TF-IDF	0.3411	0.4028	0.3694	0.3098
TX	0.3224	0.3709	0.3450	0.3147
BTM	0.3210	0.3662	0.3369	0.2927
GSDPMM	0.3440	0.4038	0.3715	0.3202
SW-IDF	0.3800	0.4460	0.4104	0.3518

Table 4. Comparison of different methods on video tag extraction of the top 5 and top 15 candidate tags

Method	H-Top 5		H-Top 15		L-Top 5		L-top 15	
	Prec	Recall	Prec	Recall	Prec	Recall	Prec	Recall
TF-IDF	0.418	0.448	0.187	0.600	0.414	0.243	0.299	0.526
TX	0.301	0.323	0.181	0.583	0.384	0.225	0.281	0.507
BTM	0.272	0.292	0.177	0.569	0.368	0.216	0.261	0.460
GSDPMM	0.281	0.301	0.183	0.593	0.418	0.249	0.307	0.539
SW-IDF	0.488	0.523	0.223	0.718	0.464	0.272	0.349	0.615

For each method, we calculate the precision, recall, mAP (Mean Average Precision) and F-score of top 10 tagging results. Results of high density and low density of comments are shown in Table 2 and Table 3 respectively.

The results show that our algorithm performs better both in high-density and low-density conditions. In high-density condition, both precision and recall are significantly increased because the semantic graph is dense. Meanwhile, noises of comments also increase with the increasing of the density of comments. Therefore the result of topic model based methods, BTM and GSDPMM are poor and even worse than classical method TF-IDF. Relatively, in low-density comments, the graph is sparse and noises reduce. Therefore the increase rate of precision and recall of our model reduce slightly. Anyhow, our performance is the highest among the others.

To further validate our algorithm, we show the precision and recall of top 5 and top 15 candidate tags in Table 4. The results of each method are similar to the performance of Top 10, which prove that our model has the best performance when extracting video tags from time-sync comments in any situation.

5. CONCLUSION

In this paper, we proposed a novel video tag extraction algorithm to acquire video tags for time-sync videos. To deal with

the features of time-sync comments, SW-IDF was designed to cluster comments into semantic association graph (SAG) by taking advantage of their semantic similarities and timestamps. In this way, the noises could be differentiated from the meaningful comments, and thus be effectively eliminated. Then, video tags were well recognized and extracted in an unsupervised way. Extensive Experiments on real-world dataset proved that our algorithm could effectively extract video tags with a significant improvement of precision and recall compared with several baselines, which obviously validated the potential of our algorithm on tag extraction, as well as tackling with the features of time-sync comments.

6. ACKNOWLEDGMENTS

This work was supported by National Natural Science Foundation of China 61532013 and 61572262; National China 973 Project No.2015CB352401; Shanghai Scientific Innovation Act of STCSM No.15JC1402400 and 985 Project of Shanghai Jiao Tong University with No.WF220103001.

7. REFERENCES

- [1] Stefan Siersdorfer, Jose San Pedro, and Mark Sanderson, "Automatic video tagging using content redundancy," in *the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2009, pp. 395–402.
- [2] Adrian Ulges, Christian Schulze, Markus Koch, and Thomas M Breuel, "Learning automatic concept detectors from online video," *Computer Vision and Image Understanding*, vol. 114, no. 4, pp. 429–438, 2010.
- [3] Bin Wu, Erheng Zhong, Ben Tan, Andrew Horner, and Qiang Yang, "Crowdsourced time-sync video tagging using temporal and personalized topic modeling," in *the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 2014, pp. 721–730.
- [4] Xiulong Liu Song Guo Liqiu Gu, Kun Wang and Bo Liu, "A reliable task assignment strategy for spatial crowdsourcing in big data environment," in *IEEE ICC 2017, Paris, France*. IEEE, 2017.
- [5] Song Guo Hongbin Chen V. C. M. Leung Kun Wang, Liqiu Gu and Yanfei Sun, "Crowdsourcing-based content-centric network: a social perspective," in *IEEE Network*. IEEE, 2017.
- [6] Santo Fortunato, "Community detection in graphs," *Physics reports*, vol. 486, no. 3, pp. 75–174, 2010.
- [7] Yikun Xian, Jiangfeng Li, Chenxi Zhang, and Zhenyu Liao, "Video highlight shot extraction with time-sync comment," in *the 7th International Workshop on Hot Topics in Planet-scale mOBile computing and online Social neTworking*. ACM, 2015, pp. 31–36.
- [8] Guangyi Lv, Tong Xu, Enhong Chen, Qi Liu, and Yi Zheng, "Reading the videos: Temporal labeling for crowdsourced time-sync videos based on semantic embedding," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [9] Zechen Wu and Eisuke Ito, "Correlation analysis between user's emotional comments and popularity measures," in *Advanced Applied Informatics (IIAIAI), 2014 IIAI 3rd International Conference on*. IEEE, 2014, pp. 280–283.
- [10] Ming He, Yong Ge, Le Wu, Enhong Chen, and Chang Tan, "Predicting the popularity of danmu-enabled videos: A multi-factor view," in *International Conference on Database Systems for Advanced Applications*. Springer, 2016, pp. 351–366.
- [11] Rada Mihalcea and Paul Tarau, "TextRank: Bringing order into texts," Association for Computational Linguistics, 2004.
- [12] David M Blei, Andrew Y Ng, and Michael I Jordan, "Latent dirichlet allocation," *Journal of Machine Learning Research*, vol. 3, pp. 993–1022, 2003.
- [13] Xiaohui Yan, Jiafeng Guo, Yanyan Lan, and Xueqi Cheng, "A biterm topic model for short texts," in *the 22nd international conference on World Wide Web*. ACM, 2013, pp. 1445–1456.
- [14] Jianhua Yin and Jianyong Wang, "A dirichlet multinomial mixture model-based approach for short text clustering," in *the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014, pp. 233–242.
- [15] Jianhua Yin and Jianyong Wang, "A model-based approach for text clustering with outlier detection," in *Data Engineering (ICDE), 2016 IEEE 32nd International Conference on*. IEEE, 2016, pp. 625–636.
- [16] Yuhua Li, David McLean, Zuhair A Bandar, James D O'shea, and Keeley Crockett, "Sentence similarity based on semantic nets and corpus statistics," *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 8, pp. 1138–1150, 2006.
- [17] Richard Socher, Cliff C Lin, Chris Manning, and Andrew Y Ng, "Parsing natural scenes and natural language with recursive neural networks," in *the 28th International Conference on Machine Learning (ICML-11)*, 2011, pp. 129–136.
- [18] Matt J Kusner, Yu Sun, Nicholas I Kolkin, and Kilian Q Weinberger, "From word embeddings to document distances," in *the 32nd International Conference on Machine Learning (ICML 2015)*, 2015, pp. 957–966.
- [19] Tomas Mikolov, Quoc V Le, and Ilya Sutskever, "Exploiting similarities among languages for machine translation," *arXiv preprint arXiv:1309.4168*, 2013.
- [20] Robert Endre Tarjan, "Efficiency of a good but not linear set union algorithm," *Journal of the ACM (JACM)*, vol. 22, no. 2, pp. 215–225, 1975.