

# From Topics to Opinions: Modelling Subjectivity for Retweeting Analysis on Twitter

## Abstract

In this paper, we investigate how user's subjectivity influence their information diffusion behavior. Inspired by psychological research, we define a general subjectivity model by combining both topics and opinions articulated in User-Generated Content (UGC) and propose an efficient framework to establish the subjectivity model. We also put forward a new way to measure subjectivity similarity between two subjectivity models. For the retweeting behavior analysis, three situations (attractiveness, sociality and popularity) are considered based on the subjectivity similarities among a target tweet, its author and followers. In the experiments, we demonstrate a user is more likely to retweet a message considering the influence of the three subjectivity similarities and the utility of our model in retweeting analysis is verified qualitatively and quantitatively on real Twitter dataset.

## Introduction

Information diffusion has drawn considerable research attentions from computer scientists, especially in the area of online social networks. Researchers have built standard models for the general information diffusion, which are useful for simulating the information flow (Goldenberg, Libai, and Muller 2001; Kempe, Kleinberg, and Tardos 2003), or detecting the outbreak of information cascades (Cheng et al. 2014). In this work, we target at a different problem: given a new message, we intend to predict which users will participate in the future diffusion process of this message (Bian, Yang, and Chua 2014). An illustration of the problem in a heterogeneous social network of Twitter can be found in Figure 1. In this example, the users have tweeted about two topics: cellphone "Iphone" and movie "Frozen". Now Tony posts a new tweet about movie "Frozen", we want to find out which one is more likely to disseminate it along all the receivers of the new tweet.

As the participants of information diffusion, humans naturally make communication and interaction by expressing opinions and preferences about the topics that interests them. In psychology, it has been identified that the subjective initiative nature of human determines that subjectivity will undoubtedly influence human's behaviors (Moore and Haggard 2008). According to theory of Biased Assimilation, people

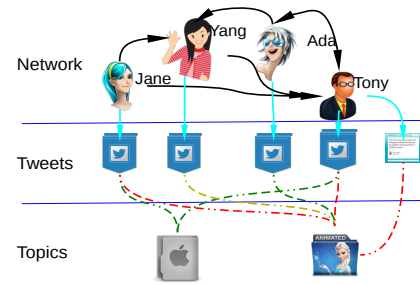


Figure 1: Motivating example. For opinions of different users, the color “red” stands for positive evaluation, “green” for negative, and “yellow” for neutral.

tend to choose and disseminate information according to their own biased opinions (Hyman 2000). Therefore, opinion and sentiment comprehension are a key aspect of users interaction in the process of information diffusion.

The propagation behaviors are different in different social networks. Twitter plays an important role in the process of information dissemination on the Internet because the retweeting convention provides an unprecedented mechanism for the spread of information despite the restricted length of a single message. Actually almost a quarter of the tweets are retweeted from others (Yang et al. 2010). Understanding how retweeting behavior works can help explaining information diffusion on Twitter.

Previous studies have developed a variety of techniques and models to capture the factors of retweeting behaviors (Macskassy and Michelson 2011; Feng and Wang 2013). However few studies have investigated the subjective motivation of a user to retweet a message. From the point of a user, retweeting is a process that includes following, evaluating and deciding whether to share. The crucial part is to evaluate whether a tweet contains information worthy enough to be shared. Therefore modelling the subjective motivation of users will provide an important perspective for retweeting behavior analysis. Intuitively, based on the principle of “like attracts like”, a biased user is more prone to retweet a message that meets his own tastes. In Figure 1, the tweets of the users present their different opinions about two topics. Tony and Jane were positive about movie “Frozen”, while Ada

was negative and Yang was neutral. For the new tweet of Tony which is positive about “Frozen”, Jane is more likely to retweet it because they both like the movie.

For the problem to be investigated, there are two questions arising: how to accurately model the subjectivity of users in terms of topics and opinions, and how to effectively measure the worthiness for the users to retweet? Answering the questions is non-trivial. In this paper, we propose a general method to model subjectivity of users, define a novel similarity measurement to calculate the worthiness, identify factors that influence a user’s retweeting behavior considering his subjectivity.

The rest of the paper is organized as follows: firstly related works are described; we give the definition and establishment details of the proposed subjectivity model before the subjectivity similarity is defined; then the factors are specified for the retweeting analysis problem; following are experiments of quantitative evaluation; and we summarize the paper and points out future work finally.

## Related Work

A large body of studies have analyzed characteristics of retweeting behavior(Bian, Yang, and Chua 2014; Luo et al. 2013), examining factors that lead to increased retweetability(Suh et al. 2010; Comarella et al. 2012) and designing models to estimate the probability of being retweeted(Jenders, Kasneci, and Naumann 2013; Pfitzner, Garas, and Schweitzer 2012). However, all of the above works neglect the subjectivity of users, which is the underlying reason for the retweeting behaviors.

Previous researches of sentiment analysis have mainly focused on reviews(Liu 2012). Recently, there have been many works on sentiment analysis for informal social media language, mainly focusing on the message level(Jiang et al. 2011; Tan et al. 2011; Guerra, Meira, and Cardie 2014). Topic models can also be utilized in sentiment analysis to correlate sentiment with topics. For example, Mei et al. (Mei et al. 2007) and Lin et al. (Lin and He 2009) attempted to incorporate the sentiment factor into topic models. Usually they learn a general word-sentiment distribution to model the sentiment of blogs or reviews, which may not work well for short and informal social media languages. Sentiment expression is deemed to be more challenging as sentiment is often embodied subtle linguistic mechanisms such as: negation, capitalization, repeated letters, exclamation and emoticon(e.g. “happy!!”), intensifiers (e.g. “like” versus “like very much”) and diminishers (e.g. “excellent” versus “rather excellent”), etc.(Brody and Diakopoulos 2011). These are hard to be modeled with probabilistic distribution. However, rule-based sentiment analysis methods can catch such subtle sentiment expressions by transforming them into rules(Thelwall et al. 2010). In our work, we adopt a rule-based method for sentiment analysis.

## Subjectivity Model

Subjectivity has been extensively studied by psychologists to characterize the personality of a person based on his historical behaviors and remarks(Engbert et al. 2007). Lin-

guists define the subjectivity of language as speakers always show their perspectives, attitudes and sentiments to events, people, topics, and entities in their linguistic contents(Stein and Wright 2005). With the explosion of social media over the past decade, more and more User-Generated Content (UGC) is available on the Web containing users’ opinions. In the Natural Language Processing area, opinion mining techniques(Liu 2012) have been developed to computationally model the subjectivity of users. A variety of aspect-based or topic-sentiment models have been built from UGC by casting opinions as polarity, ratings, or emotions regarding a topic(Lek and Poo 2013; Mei et al. 2007). But they are often limited in utility by their definition of opinions. We give a general framework to model subjectivity by combining topics and opinions together with a new representation of opinions. Here we give our definition of subjectivity model under context of Twitter, while we emphasize that our model can be adapted to other platforms as well.

## Definition

Let  $G = (V, E)$  denote a social network on Twitter, where  $V$  is a set of users, and  $E \subset V \times V$  is a set of follow relationships between users. For each user  $u \in V$ , there is a tweets collection  $M_u$  denoting his message history. We assume that there is a topic space  $T$  containing all topics users in  $V$  talk about, and a sentiment space  $S$  to evaluate their opinions towards these topics. For the “subjectivity” of a user  $u \in V$ , we refer to both topics and opinions articulated in his tweets collection  $M_u$ .

**Definition 1 (Subjectivity Model)** *The subjectivity model of user  $u$ , is the combination of topics of interest  $\{k\}$  in topic space  $T$  and his opinions  $\{O_k\}$  towards each topic distributed over sentiment space  $S$ .*

$$SM_u = \{(k, w_{u,k}, \{d_{u,k,s} | s \in S\}) | k \in T\}$$

where:

- with respect to user  $u$ , for each topic  $k \in T$ , weight  $w_{u,k}$  represents the distribution of the user’s interests on it, subject to  $\sum_{k=1}^{|T|} w_{u,k} = 1$ .
- opinion of  $u$  towards topic  $k$  is modeled as a topic-dependent sentiment distribution over sentiment space  $S$ ,  $O_k = \{d_{u,k,s} | s \in S\}$ , subject to  $\sum_{s=1}^{|S|} d_{u,k,s} = 1$ .

Our model is more general than others in that we combine the topics of interest and topic dependent opinions into a holistic framework, and more importantly, we define opinion as a probabilistic distribution in a scalable sentiment space. The sentiment space can cover all the sentiment modalities. For example, it could be binary space standing for sentiment polarity, or sequential space for sentiment strength, or discrete space for emotions. Figure 2 is a visualized subjectivity model example in a  $[0, 100]$  topic space and a  $[0, 8]$  sentiment space.

## Establishment of Subjectivity Model

In this section, we present our method to concrete subjectivity model by deriving topics and opinions from the UGC of all users  $M = \{M_u | u \in V\}$ .

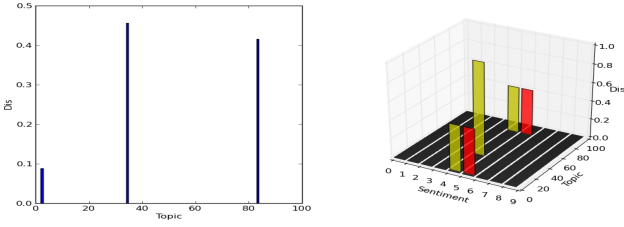


Figure 2: Subjectivity model example. The left subgraph denotes interests distribution on topic 2, 32 and 83: ( $w_{u,2} = 0.08, w_{u,32} = 0.48, w_{u,83} = 0.44$ ). The right subgraph denotes opinions towards topics:  $O_2 = (d_{u,2,4} = 0.5, d_{u,2,5} = 0.5)$ ,  $O_{32} = (d_{u,32,4} = 1.0)$ ,  $O_{83} = (d_{u,83,4} = 0.5, d_{u,83,5} = 0.5)$ .

**Topic Analysis** We simply use the concept of topics to broadly refer to the different kinds of content such as key words, entities, hashtag, etc. Previous studies have tried to identify topics from tweets by finding key words(Chen et al. 2010), extracting entities(Abel et al. 2011) or linking tweets to external knowledge categories(Macskassy and Michelson 2011). However, works show that topic model is more effective in identifying topics from short and informal social media language(Hong and Davison 2010). In this paper, state-of-the-art Latent Dirichlet Allocation (LDA)(Blei, Ng, and Jordan 2003) is employed for unsupervised topic discovery and for topic assignment of future tweets. LDA can be used to find a set of  $K$  latent topics from a document corpus, and then to represent each document  $D$  with a distribution  $\theta_D$  of the latent topics. For each word  $w_i$  in  $D$ , a topic  $z$  is first sampled from the document topic distribution  $\theta_D$ , then  $w_i$  is sampled according to word distribution  $\phi_z$  of topic  $z$ .

As the first step, we adopt the user-level LDA model to build a global Topic Model (TM), which regards all tweets of a user as one document of LDA(Hong and Davison 2010). The TM model will be used throughout our framework. Usually a tweet concentrates a single topic within its short length, therefore we assign a tweet  $t$  to a topic that maximizes the probability of generating  $t$ :

$$z_t = \arg \max_k \prod_{w \in t} P(w|\phi_k) \quad (1)$$

We can get the weight distribution on each topic  $k$  of user  $u$  by normalizing all tweets that talk about topic  $k$ :

$$w_{u,k} = \frac{|\{t : t \in M_u \wedge z_t = k\}|}{|M_u|} \quad (2)$$

**Opinion Analysis** In Figure 1, both Tony and Jane held overall positive opinion about the movie “Frozen”, but maybe they liked the movie for different reasons. Jane mainly liked the romantic story of this movie but was a little disappointed about its animation picture, while Tony liked this movie because he was mostly convinced by its animation technology although he disliked the prince and princess genre. If we represent their opinions with a simple binary polarity, without differentiating their opinions on different

aspects, the subjectivity model may not satisfy the information discovery needs of different users. Therefore, it is better to describe opinion for a topic as a probability distribution over the sentiment space. Furthermore, a more fine-grained sentiment space is preferred if we want to distinguish subjectivities of users more precisely.

Researches on sentiment analysis of social media have provided many effective state-of-the-art techniques and tools(Thelwall et al. 2010; Hu et al. 2013), with which sentiment of a tweet  $t$  can be identified as  $s_t$ . The opinion distribution  $O_k$  toward a topic  $k$  is:

$$\begin{aligned} O_k &= \{d_{u,k,s} | s \in S\} \\ &= \left\{ \frac{|\{t : t \in M_u \wedge z_t = k \wedge s_t = s\}|}{|M_u|} | s \in S \right\} \quad (3) \end{aligned}$$

### Subjectivity Similarity

With the subjectivity model established, a subjectivity similarity measurement needs to be calculated to analyze various subjective decision-making process such as retweet. Firstly we should define the opinion similarity on a common topic.

**Opinion Similarity** Opinion in the subjectivity model is treated as a distribution over sentiment space with each dimension of the distribution representing the proportion of the corresponding sentiment value. In fact, values of the sentiment space are not independent. They are sequential in magnitude and quantized to measure the strength of sentiment. Therefore, normal distribution similarity measurements such as KL-divergence and cosine similarity are not suitable for such kind of opinion distribution. As illustrated in Table 1, in a  $[0, 8]$  integer sentiment space, opinion  $O_k^1$  is most negative (100% of value 0), opinion  $O_k^2$  (100% of value 6) is positive, and  $O_k^3$  (100% of value 8) is most positive. If the cosine similarity measurement is adopted to calculate opinion similarity, all similarities among them are 0. In fact  $O_k^2$  is more similar with  $O_k^3$  than  $O_k^1$  because they both hold positive opinion and their sentiment strength distance is much less than with  $O_k^1$ . Therefore, opinion similarity can’t be calculated simply as the normal probabilistic distributions, or just as the sentiment strength distance. To accurately catch opinion similarity, we propose a novel method by combining both sentiment distance and distribution similarity. The opinion similarity between two opinions  $O_k^1, O_k^2$  on the same topic  $k$  can be calculated as:

$$Sim(O_k^u, O_k^v) = \frac{|S| - |\sum_{i=0}^{|S|} d_i^u v_i - \sum_{i=0}^{|S|} d_i^v v_i|}{|S|} \quad (4)$$

where  $d_i$  denotes the  $i^{th}$  dimension of opinion distribution, and  $v_i$  denotes corresponding sentiment strength value. The

Table 1: Illustration of opinion similarity

	0	1	2	3	4	5	6	7	8
$O_k^1$	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
$O_k^2$	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0
$O_k^3$	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0

similarities of opinions in Table 1 calculated with Equation 4 are  $Sim(O_k^1, O_k^3) = 0$ ,  $Sim(O_k^2, O_k^3) = 6/8$  and

$Sim(O_k^1, O_k^2) = 2/8$ , which are consistent with our intuitive understanding.

**Subjectivity Similarity** As the subjectivity model indicates, a user's topics of interest is represented with a weight distribution over topic space  $T$ . Therefore, the subjectivity similarity should be integrated from the weight distribution and the opinion similarities on each topic:

$$Sim(SM_u, SM_v) = \frac{\sum_{k=1}^{|T_{u,v}|} \theta_u(k) Sim(O_k^u, O_k^v)}{\sum_{k=1}^{|T_{u,v}|} \theta_u(k)} \quad (5)$$

Where  $T_{u,v}$  denotes the common topics between two users  $u$  and  $v$ , which can be regarded as the intersection of their topics of interest;  $\theta_u(k)$  denotes the weight of topic  $k$ , and  $\sum_{k=1}^{|T_{u,v}|} \theta_u(k)$  is the normalized factor.

Note that, when we measure how similar user  $u$  is with user  $v$ , we use the topic weights of user  $u$ , thus the subjectivity similarity is asymmetric. The intuition lies in that subjectivity of a user is a personal inner interest taste, and it is also subjective judgement about how like-minded a friend is. Therefore, for the measurement of subjectivity similarity,  $Sim(SM_u, SM_v) \neq Sim(SM_v, SM_u)$ .

## Retweeting Analysis

Whether a user propagate a tweet may be affected by various factors. From the point of a user, three situations usually make him retweet: 1) the content of the tweet is attractive for the user, and his retweeting behavior is in accordance with subjectivity; 2) the tweet is posted by the user's close friend, and his retweeting behavior is due to social needs; and 3) the content is novel or epidemic, and his retweeting behavior is a result of conformity needs (i.e., the act of matching opinions to group norms (Cialdini and Goldstein 2004)). These situations exhibit different types of reasons a user retweet a message, and we will analyze them by quantifying them as three subjectivity similarities in this section. For a target tweet  $t$ , let  $F$  denote the followers who receive  $t$  by following its author  $u_a$ . For each follower  $f \in F$ , we can define a quadruple  $\langle f, u_a, m, r_f \rangle$ ,  $r_f$  is a binary label indicating if  $t$  is retweeted by  $f$ , which needs to be predicted based on our analysis.

### Attractiveness

A user is likely to repost a microblog if the user finds the content is attractive according to his subjective judgement. We can measure such attractiveness quantitatively by calculating the subjectivity similarity between the tweet  $t$  and user  $f$ . For a tweet  $t$ , its topic  $z_t$  can be identified with Equation 1, and let  $s_t$  be its sentiment. The content of  $t$  can also be modeled using our subjectivity model definition with a single topic of interest and a 100% opinion distribution on a single sentiment value. Thus the attractiveness of the tweet can be measured with subjectivity similarity using Equation 5, which is marked as  $Sim(f, t)$ .

### Sociality

In this case, the retweeting behavior is based on the needs of social interaction. That is, the behavior is triggered because

the tweet  $t$  is sent by a like-minded friend, instead of the information it contains. We can measure how like-minded the user  $f$  and his friend  $u_a$  with their subjectivity similarity  $Sim(f, u_a)$ . However, different kinds of friends may have a different influence on the user  $f$ . For example,  $f$  may follow many friends, but only frequently interacts with a few. Furthermore, not all tweets of a friend may be of interest to  $f$ . For example, in Figure 1, Jane may be interested in tweets from Tony about movie, but may not be interested in his tweets about cellphone. We therefore assign a weight to  $Sim(f, u_a)$  to reflect the influence of different kinds of friend, which is composed of four factors:

**Expert Factor**  $w_E(u_a)$  It represents the relative expertise of the author  $u_a$  among his followers including  $f$ . The expert user imposes more influence on others. We simply calculate it as the ratio of user  $u_a$ 's tweets count over all tweets of  $u_a$  and his followers by  $w_E(u_a) = |M_{u_a}| / |\{M_u | u \in u_a \cup F\}|$ .

**Leadership Factor**  $w_L(u_a)$  In our work, the leadership of a user  $u_a$  is determined by his followers. The leadership weight is calculated by  $w_L(u_a) = \log(|F|) / \log(\max)$ , where  $\max$  is the maximum popularity of a user in Twitter<sup>1</sup>.

**Similarity Factor**  $w_S(u_a, f)$  The similarity of interests between  $u_a$  and  $f$  is measured as the inverse KL-divergence between their topic distribution in their subjectivity model:  $w_S(u_a, f) = 1 / KL(\theta_{u_a}, \theta_f)$ .

**Interaction Factor**  $w_I(u_a, f)$  All the interactions  $Interaction_{u_a, f}$  between  $u_a$  and  $f$  are analyzed, which include the conversations between them, mentions of each other, and retweets from each other. The factor weight is calculated by normalizing  $Interaction_{u_a, f}$  with all tweets of  $u_a$  and  $f$ :  $w_I(u_a, f) = |Interaction_{u_a, f}| / |\{M_{u_a}, M_f\}|$ .

Above all, the influence weight is the combination of four factors:

$$w_{u_a, f} = \lambda_1 * w_E(u_a) + \lambda_2 * w_L(u_a) + \lambda_3 * w_S(u_a, f) + \lambda_4 * w_I(u_a, f). \quad (6)$$

where  $\lambda_i$  is an optional weight vector to enable different influence of the factors, subject to  $\sum_{i=1}^4 \lambda_i = 1$ . We set them uniformly as 0.25 in the experiment.

### Popularity

If a tweet is popular (e.g., breaking news), it will be very probable to be retweeted. In this situation, the tweet  $t$  is often inconsistent with the interests and opinion of its author  $u_a$ . Thus the similarity between  $t$  and  $u_a$  in terms of subjectivity is very low, which is marked as  $Sim(u_a, t)$ . The retweeting behavior is highly related to the popularity of  $t$  rather than the content or the friend who post it. We assign a popularity weight  $w_P$  to  $Sim(u_a, t)$ , which is the proportion of user  $f$ 's followers who have retweeted the tweet  $t$ .

From the point of motivation, a user might retweet a message if its content is approximate to his subjectivity, its author is a like-minded friend and it is popular among his friends. In next section we carry out a set of experiments to inspect and verify the impact of such motivation on retweeting behavior.

<sup>1</sup><http://twittercounter.com/pages/100>



## Experiments

### Dataset and Settings

We adopt the Twitter dataset of a previous work (Luo et al. 2013). To form the dataset, 500 target English tweets were monitored to find who would retweet it in the future. Each target tweet was set as starting point to collect recent tweets for its author and followers. Overall, there are 45,531 users who have posted at least 6,277,736 tweets. 5214 users have retweeted at least one target tweet during the monitored period. To avoid the bias introduced by dataset imbalance, an evaluation dataset is constructed by taking 5,214 retweeters as positive instances, and randomly sampling 5,214 non-retweeters as negative instances.

For the topic model, we use ~~package~~ Gensim (Řehůřek and Sojka 2010), which adopts an efficient batch-based online inference algorithm. All parameters are set as defaults and the number of topic traverses from 50 to 200. For sentiment analysis, we just make use of an off-the-shelf work, i.e. SentiStrength (Thelwall et al. 2010). In order to catch the sentiment of tweets, we use the sentiment lexicon created based on AFINN by Nielsen (Mohammad, Kiritchenko, and Zhu 2013). The sentiment space is formed by mapping the positive and negative sentiment values to range  $[0, 8]$ .

### Correlation Test

First of all we assess the correlation between subjectivity similarity and retweeting behavior with a statistical hypothesis test Analysis of Variance (ANOVA) (Fisher et al. 1970), which tests the *null hypothesis* that the retweeters and non-retweeters have the same subjectivity similarity means. The results are listed in Table 2. The bold-faced entries mean that the *p-value* is lower than significance level. Note that

Table 2: ANOVA results for subjectivity similarities. If the difference is due to chance, *F-ratio*=1.00, otherwise *F-ratio* > 1.00 (*p-value* < 0.01).

Similarity		$Sim(f, t)$	$Sim(f, u_a)$	$Sim(u_a, t)$
50	<i>F</i>	<b>12.182</b>	2.212	4.236
	<i>p</i>	<b>4.44e<sup>-06</sup></b>	0.140	0.272
100	<i>F</i>	<b>43.892</b>	<b>31.145</b>	<b>28.466</b>
	<i>p</i>	<b>8.65e<sup>-11</sup></b>	<b>3.55e<sup>-08</sup></b>	<b>1.32e<sup>-09</sup></b>
150	<i>F</i>	<b>22.356</b>	<b>12.240</b>	<b>14.664</b>
	<i>p</i>	<b>2.43e<sup>-08</sup></b>	<b>6.25e<sup>-06</sup></b>	<b>8.46e<sup>-07</sup></b>
200	<i>F</i>	<b>31.675</b>	<b>20.616</b>	6.145
	<i>p</i>	<b>4.22e<sup>-06</sup></b>	<b>2.92e<sup>-05</sup></b>	0.26

for the topic numbers of 100 and 150, all similarities yield *p-values* below significance level with *F-ratio* above 1.00. This suggests that the subjectivity similarities could be useful features for modeling retweeting behavior. For the rest experiments, we set the topic number as 100 for LDA.

### Case Study

In this section, we give an ~~vivid~~ example to illustrate the subjectivity model and its ability in explaining the retweeting behavior. The subjectivity models of one of the 500 target tweets, its author, and two followers (one retweeter, the other non-retweeter) are shown in Figure 3. The right part of

each sub-figure illustrates topic distribution and the left part illustrates opinions towards each topic.

The tweet is about the 14<sup>th</sup> topic, and the opinion is neutral. The author concentrates on the 14<sup>th</sup> topic, and his opinion is mainly neutral. As for two followers, the retweeter has tweeted about two topics (the 14<sup>th</sup> and 52<sup>nd</sup> topic) uniformly and his opinion towards the 14<sup>th</sup> topic is mainly neutral. While the non-retweeter has also talked about two topics (14<sup>th</sup> and 56<sup>th</sup> topic), but he is mainly interested in the 14<sup>th</sup> topic and his opinion is positive.

Table 3 shows the three subjectivity similarities for both retweeter and non-retweeter. They have common interest

Table 3: Illustration of example subjectivity similarities

Similarity	$Sim(f, t)$	$Sim(f, u_a)$	$Sim(u_a, t)$
Retweeter	0.854	0.967	0.886
Non-retweeter	0.805	0.919	0.886

(the 14<sup>th</sup> topic), and furthermore the non-retweeter is more similar with the tweet and its author than the retweeter in terms of topics. But their different opinions towards the topic elicit their different behaviors, which verifies our model can help better understanding the retweeting behavior not only from topics but also opinions.

### Performance Evaluation

We carry out the retweeting classification experiments in three stages. Firstly we compare our model against other topic-based models including TF-IDF model (modeling user interests using bag-of-words), entity-based model (using entities extracted from the UGC) and hashtag-based model (using hashtags used in the UGC) (Abel et al. 2011). Secondly, our model is compared with two generative topic-sentiment models (TSM model (Luo et al. 2007) and JST model (Lin and He 2009)). TSM and JST can also model topic and topic related sentiment simultaneously. We use Equation 5 to calculate three subjectivity similarities for TSM and JST as our model, and combine them together in the classification.

Our model can catch the subjective motivation of users based on UGC, whereas other important factors associated with retweeting behavior are not considered, such as network topology and meta-data of users. Therefore, our model is also compared with the method of Luo et al. (2013) (marked as “LUO”), in which different factors that might affect retweeting behaviors are considered. They only use bag-of-words to model user interests, so we also carry out combining experiments to demonstrate that performance of prediction can be improved by replacing their bag-of-words model with our model (marks with “LUO+” prefix).

The logistic regression classifier is used for training and testing in a 5-fold cross-validation manner. We set a baseline, which simply predicts users who have retweeted the author previously as the retweeters of target tweet. All results are presented in Table 4 in terms of accuracy.

Firstly, all models except the hashtag-based model outperform the baseline (60.85%) significantly. While for hashtag-based model, the accuracy is only 59.12%, the reason lies in

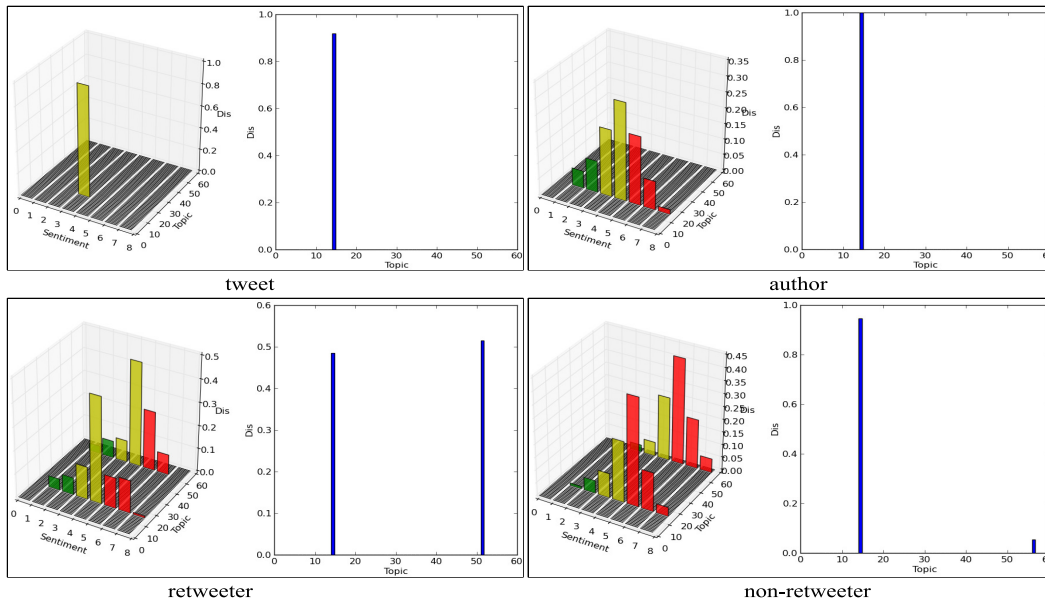


Figure 3: An illustration of subjectivity models of a tweet, author and two followers.

Table 4: Accuracy performance. A significant improvement over baseline with \* and LUO’s model with ‡ ( $p < 0.05$ ).

Feature	Accuracy(%)	Feature	Accuracy(%)
baseline	60.85		
TF-IDF	62.85 *	LUO	71.76 *
entity	68.76 *	LUO+entity	72.15 *
hashtag	59.12	LUO+hashtag	68.44 *
TSM	67.44 *	LUO+TSM	68.23 *
JST	68.13 *	LUO+JST	70.53 *
$Sim(f, t)$	73.88 * ‡	LUO+ $Sim(f, t)$	74.04 * ‡
$Sim(f, u_a)$	70.04 *	LUO+ $Sim(f, u_a)$	70.27 *
$Sim(u_a, t)$	69.64 *	LUO+ $Sim(u_a, t)$	71.86 *
$sim_{all}$	75.64 * ‡	LUO+ $sim_{all}$	78.15 * ‡

the spacity of hashtag in tweets.

Secondly,  $Sim(f, t)$  and  $sim_{all}$  outperform “LUO” (71.76%) significantly. The best performance is achieved by the  $sim_{all}$  (75.64%), for which we add three similarities to the classifier to test the impact of their combination. The performence of TF-IDF model (62.85%) is a little better than baseline. The entity-based model (68.76%) is very close to  $Sim(f, u_a)$  (70.04%) and  $Sim(u_a, t)$  (69.64%), and the difference is not significant.

Thirdly, the performance of two topic-sentiment models (TSM: 67.44%, JST: 68.13%) is not as good as our models. The reason lies in that they use a binary sentiment representation (positive or negative), which can not differentiate opinions elaborately. Our model can capture more subtle and fine-grain sentiment, which could distinguish different subective motivation of retweeting behavior.

Finally, in the combining evaluation,  $Sim(f, t)$  gives a significant improvement (LUO+ $Sim(f, t)$ , 2.12% improvement) over “LUO”, but other two similarities and

the entity-based model can not improve performance significantly. The performance is even degraded after combining with the hashtag-based model and two topic-sentiment models. But noticing that, the most significant improvement(LUO+ $sim_{all}$ , 6.39% improvement) is achieved by combining with all three similarities.

Above all, the results show that our model can better help predicting retweeting behavior and can be regarded as a useful way to analyze the retweeting behaviors of users.

## Conclusion

Motivated by the psychological research, this paper postulates that the online behaviors of social media users are affected by their subjectivity. Therefore, a general subjectivity model has been proposed and an efficient framework has been designed to establish the subjectivity model. Also a novel method is put forward to measure the subjectivity similarity. The subjectivity model has been applied to the retweeting analysis considering three different situations, and these situation are quantified with subjectivity similarities. Experiment results demonstrate the effectiveness of the proposed model in the retweeting analysis problem and show that the model is able to reach better understanding of retweeting behavior. In the future, we will apply our model to other social network analysis task such as link prediction and recommendation.

## References

- Abel, F.; Gao, Q.; Houben, G.-J.; and Tao, K. 2011. Analyzing user modeling on twitter for personalized news recommendations. In *UMAP*. Springer. 1–12.
- Bian, J.; Yang, Y.; and Chua, T.-S. 2014. Predicting trending messages and diffusion participants in microblogging net-

