# Topic Related Opinion Integration for Users of Social Media

Songxian Xie, Jintao Tang, and Ting Wang

School of Computer Science, National University of Defense Technology, Hunan, PRC
{xsongx,tangjintao,tingwang}@nudt.edu.cn

**Abstract.** Social media such as Twitter, has become a valuable source for mining opinions of users about all kinds of topics. In this paper, we investigate how to automatically integrate topic related opinions expressed by a user in User-Generated Content (UGC). We propose a general subjectivity model by combining topics and fine-grained opinions towards each topic, and design an efficient algorithm to establish the model. We demonstrate utility of our model in the opinion prediction problem and verify the effectiveness of our model qualitatively and quantitatively in a series of experiments on real Twitter data. Results show that the proposed model is effective and can generate consistent integrated opinion summaries for users. Furthermore, the proposed model is more suitable for social media context, thus can reach better performance in an opinion prediction task.

**Keywords:** LDA, social media, opinion integration, subjectivity model.

## 1 Introduction

With the rise of content-based social media such as Twitter, millions of users are more and more willing to publish online short messages to express their opinions on a great variety of topics they are interested in. The wide coverage of topics, dynamics of discussion, and abundance of opinions imbedded in the social media data make them extremely valuable source for mining users' opinions about all kinds of topics (e.g., products, political figures, etc.), which in turn can enable a wide range of applications, such as opinion search for ordinary users, opinion tracking for business intelligence, and user behavior prediction for targeted advertising. However, with such a large scale of information source, it is quite challenging to integrate and digest all the opinions from different users. For example, a query "iPhone" on Twitter (as of Jan. 14, 2014) returns 830,879 tweets of 231,233 users, suggesting that there are many users have expressed opinions more than once about iPhone in their tweets. To enable an application to benefit from all kinds of opinions of different users, it is thus necessary to automatically integrate and present an overall opinion summary for each user [1]. In fact, users often publish several messages on the topics they are interested, therefore how to find these topics and integrate opinions towards each topic scattered in many independent tweets of a user poses special challenge for opinion mining related researchers.

In this paper, we propose a combining model (named as subjectivity model) by incorporating topics and opinions at the user level, of which one part represents topics of interest distribution, while the other part represents the distribution of opinions towards these topics. Specifically, we propose a general method to solve this integration problem in three steps illustrated as in Figure 1: (1) extract topics of interest from tweets of a user using user-level LDA; (2) extract separate opinion and topic for each tweet with sentiment and topic analysis (3) summarize and integrate the extracted opinions towards each topic to form a subjectivity model for each user.
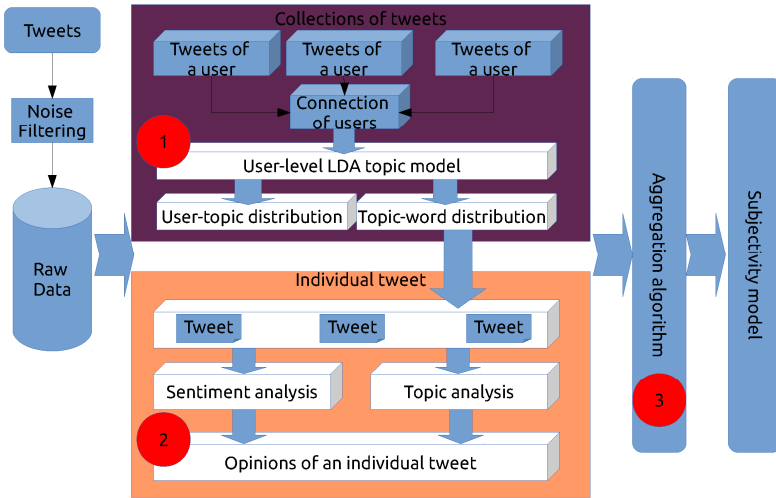


**Fig. 1.** Framework

The rest of the paper is organized as follows. In Section 2, we introduce related works, and we formally define the novel problem of opinion integration in Section 3. After that, we present our model and analyze the difference with generative model in Section 4. We discuss our experiments and results in Section 5. Finally, we conclude our work and point out future target.

## 2    Related Works

Sentiment analysis is a popular research area and previous researches have mainly focused on reviews or news comments [2, 3]. Recently, there have been many works on sentiment analysis on Twitter, mainly focusing on the tweet level [4, 5, 6, 7, 8], of which, the techniques employed are generally standard tweet-level algorithms that ignore many special characteristics of social media. There have been also some previous works on automatically determining user-level opinions or ideology [9, 10], generally looking at information embedded in the contents that the users generate. Most of related researches mainly focused on identification of sentimental object [11], or detection of objects' sentimental polarity [12] without considering the topic aspects.

Since the introduction of topic model such as LDA [13], various extended models have been used for topic extraction from large-scale corpora at user level [14, 15]. Topic models can also be utilized in sentiment analysis to correlate sentiment with topics. Mei et al. [16] and Lin et al. [17] incorporated topic models and sentiment analysis for reviews and blogs.

## 3    Opinion Integration Problem

As we describe in Section 1, a user usually posts multiple messages on various topics during his social media usage. Therefore what's the opinion of a user on a specific topic can't be determined from just one tweet, but should be integrated from all the topic related tweets he has posted. In this paper, we put forward a new problem which is defined as Opinion Integration Problem (OIP). We focus on user-level rather than tweet-level opinion because the end goal of opinion mining technologies is to find out what a person thinks but not what only a piece of message states, and the identification of the opinion articulated in an individual text is usually a middle step for that ultimate objective. Additionally, it is plausible that there are cases where opinions of a user in one tweet is ambiguous because they are restricted to be so short that the context of its opinion is missing, but his overall opinion can be determined by looking at his collection of tweets [8].We illustrate a typical scenario of user-level topic related opinion integration problem on Twitter in Figure 2.
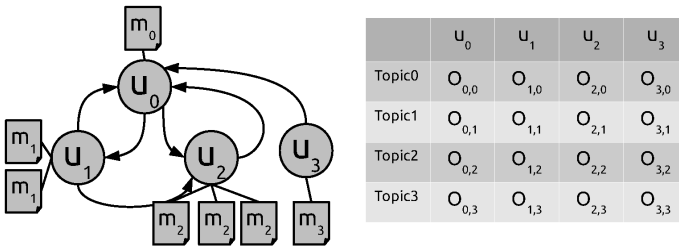


|        | $u_0$ | $u_1$ | $u_2$ | $u_3$ |
|--------|-------|-------|-------|-------|
| Topic0 | $O_{0,0}$ | $O_{1,0}$ | $O_{2,0}$ | $O_{3,0}$ |
| Topic1 | $O_{0,1}$ | $O_{1,1}$ | $O_{2,1}$ | $O_{3,1}$ |
| Topic2 | $O_{0,2}$ | $O_{1,2}$ | $O_{2,2}$ | $O_{3,2}$ |
| Topic3 | $O_{0,3}$ | $O_{1,3}$ | $O_{2,3}$ | $O_{3,3}$ |

**Fig. 2.** Illustration of Opinion Integration Problem

**Definition 1(Opinion Integration Problem).** *As shown in the figure, there is a heterogeneous network of Twitter consisted of users set $V = \{u_i\}$, directional relations set $E = \{(u_i, u_j) \mid u_i, u_j \in V\}$ of all users, and their associated tweets $M_i = \{m_i\}$, in which topics (denoted as $T = \{Topic_k\}$) and opinions of tweets can be determined and extracted. For a user $u_i$, his opinion (denoted as $O_{i,k}$) towards topic $Topic_k$ is not the opinion imbedded in his single tweet $m_i$, but the integrated opinion from all his tweets $M_i = \{m_i\}$.*

There are two important factors that must be taken into considerations for the OIP problem. Firstly, topics both users and tweets talk about should be determined in a same topic space so as the target of opinion is consistent. Secondly but most importantly, opinions and topics are closely related, tweets of a user around some topic

often cover a mixture of aspects related to that topic with different preferences. Different opinions may be expressed by the user towards different aspects, where users may like one aspect of a topic but dislike others. Therefore, how to integrate all opinions of tweets related to a topic into one holistic opinion and represent it reasonably poses special challenge. In this paper, we propose a novel subjectivity model to meet these two challenges.

# 4 Subjectivity Model

In this section, we give a formal definition of the model we work with to meet the challenges of OIP problem, which has been substantially defined and described in our previous work [18]. Here we only repeat the definition and the algorithm of model establishment, for more details, please refer to our paper [18]. Usually user level opinion is to classify each user's sentiment on a specific topic into one of two polarities: "Positive" and "Negative". "Positive" means that the user supports or likes the target topic, whereas "Negative" stands for the opposite. However in our model we adopt a broad "opinion" definition as sentiment coverage towards a topic over a fine-grained sentiment values to differentiate subtle opinions of users. For example one is more positive about a topic with sentiment strength 8 than another user with sentiment strength 7. At the same time, we define opinion of a user as a probabilistic distribution over the sentiment values instead of one single value, considering the user may express his different opinion on different aspects of the same topic. The notion of "opinion" is quite vague; we adopt this broad definition to ensure generality of the model. We frame the model in the context of Twitter to keep things concrete, although adaptation of our model to other social network settings is straightforward. We name our model as "subjectivity model" as it models the subjective information in the content generated by a user. Therefore, we give a formal definition of the subjectivity model under the context of Twitter as follows.

## 4.1 Definition

Let $G = (V, E)$ denotes a social network on Twitter, where $V$ is a set of users, and $E = V \times V$ 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 user $u \in V$, we refer to both topics and opinions articulated in his tweets collection $M_u$.

**Definition 2 (Subjectivity Model).** *The subjectivity model $P(u)$ of user $u$, is the combination of topics $\{t\}$ the user talks about in topic space $T$ and his opinions $O_t$ towards each topic distributed over sentiment space $S$.*

$$P(u) = \{(t, w_u(t), \{d_{u,t}(s) \mid s \in S\}) \mid t \in T\} \tag{1}$$

*where:*

- *with respect to user $u$, for each topic $t \in T$, its weight $w_u(t)$ represents the distribution of the user's interests on it, subject to $\sum_{t=1}^{|T|} w_u(t) = 1$.*
- *opinion of the user towards topic $t$ is modeled as a topic related sentiment distribution over sentiment space $S$, $O_t = \{d_{u,t}(s) \mid s \in S\}$, subject to $\sum_{s=1}^{|S|} d_{u,t}(s) = 1$.*

Subjectivity model aims at obtaining the topic related refined sentiment for investigating user-level opinion mining, which can get a comprehensive understanding of the subjectivity for a user by modeling both his topics of interest and opinions towards each topic.

## 4.2     Establishment of Subjectivity Model

According to the definition of subjectivity model, there are two distributions to model the subjectivity: the topic distribution and the opinion distribution for each topic. Both of them need to be inferred from historic content produced by users.

For users set $V$ of a social network, we denote tweets set published by a user $u \in V$ as $M_u = \{m_u\}$. $M_u$ is concatenated to a document $d_u$ to construct topic space $T = \{t_i \mid i = 1, \cdots K\}$ with user-level LDA model. The topic model is built with parameter $\theta$ representing the distribution of each user over topics in the topic space $T$, and parameter $\beta$ representing the distribution of each topic over the vocabulary of all tweets. SentiStrength [25] is applied to each tweet $m$ in collection $M_u$ and outputs sentiment strength $s_m$ for tweet $m$. With statistical topic analysis and opinion analysis for each user and tweet, we put forward a novel algorithm to concrete subjectivity model $P(u)$ for user $u$ as algorithm 1. In the algorithm, we assume the sentiment of tweet m is related to every topic it talks about in $Z_m$ for simplicity.

---

**Algorithm 1**. Establishment of subjectivity model.

**Input:** The users set of asocial network $V$;
    The tweets set published by each user $u$, $M_u$;
**Output:** The subjectivity model for each user $u$, $P(u)$;
Topic analysis with a user-level LDA, getting a topic model $P(\theta, \beta \mid M_u, V)$;
**for all** tweet $m \in M_u$ **do**
    Sentiment analysis, outputting sentiment of $m$, $s_m$;
**end for**
**for** user $u \in V$ **do**
    the topic distribution is the corresponding component of parameter $\theta$, $\theta_u$;
    the topics $u$ tweets about are $Z_u = \{t \mid p(t \mid \theta_u) > 0, t \in T\}$;

**end for**

**for** $m \in M_u$ **do**

  topics of $m$ can be identified by the topic model:

$$Z_m = \{t \mid p(t \mid \theta, \beta, Z_u) > 0, t \in T\} \tag{2}$$

**end for**

**for** each topic $t \in Z_u$ **do**

  **for** sentiment value $s \in S$ **do**

    count the number of tweets that talk about topic $t$ with sentiment value $s$:

$$N_s = \sum_{m \in M_u \wedge s_m = s \wedge t \in Z_m} I(s_m) \tag{3}$$

  **end for**

  calculating opinion towards topic $t$:

$$O_t = \left\{ \frac{N_s}{\sum_{s \in S} N_s} \mid s \in [0, S] \right\} \tag{4}$$

**end for**

establishing subjectivity model of user $u$:

$$P(u) = \{(t, p(t \mid \theta_u), \left\{ \frac{N_s}{\sum_{s \in S} N_s} \right\}) \mid t \in Z_u, s \in S\} \tag{5}$$

**return** $P(u)$.

### 4.3    Application of Subjectivity Model

The learned subjectivity model can be used to help with many applications such as opinion mining and behavior prediction (retweet, follow, etc.). Here we demonstrate one application on, i.e., how the learned model can help improve the performance of user opinion prediction. Our strategy is based on the premises that users usually tend to express their opinions consistently. In other words, positive and negative opinions are not randomly expressed by people. E.g., a user who supports a candidate in an election will tend to post positive tweets on a regular basis. Technically, social theories say that the user exhibits a varying degree of bias, which is his subjectivity [19].

We formulate the opinion prediction of a user as a triplet in the form of $< author, m, t >$, where author is the user who post tweet m, which talks about topic $t$. The goal is to predict the polarity $p = \{positive, negative\}$ of tweet $m$ toward topic $t$. For such a problem, the dominant approach relies on extracting textual patterns from the tweet $m$ and exploiting these patterns to predict its polarity.

However subjectivity model of a user provides information that is more robust to a single tweet short of context, as it is more consistent than typical textual information.

Thus, we propose an alternative approach to improve the performance of opinion mining of a single tweet based on subjectivity model of its author. Specifically, for tweet $m$, subjectivity model of its author $P(author)$ can be concreted according to algorithm 1. Let $s_m$ denote its sentiment value calculated with some sentiment classifier such as SentiStrength. The topic tweet m talks about can be identified with equation 2 in algorithm 1:

$$\hat{t} = \arg\max(\hat{P}(t \mid \theta, \beta, Z_u) \mid t) \tag{6}$$

Thus opinion distribution of the user author can be identified from his subjectivity model $P(author)$: $O_{author,\hat{t}}$, which is a distribution over sentiment value space $S$. We can get a normalized sentiment value of the user on topic $\hat{t}$:

$$\hat{s}_m = \sum_{i \in T} d_i * v_i \tag{7}$$

where $v_i$ denotes the sentiment value and $d_i$ denotes the corresponding dimension of the sentiment distribution. Now we can predict the polarity $p$ by smoothing the sentiment of tweet $m$ with the normalized sentiment value of its author:

$$p = \begin{cases} positive & if \quad \dfrac{\hat{s}_m + s_m}{2} > \dfrac{|S|}{2} + 1; \\ negative & if \quad \dfrac{\hat{s}_m + s_m}{2} < \dfrac{|S|}{2}; \\ neutral & otherwise. \end{cases} \tag{8}$$

## 5      Experiment

### 5.1      Dataset and Settings

We use an off-the-shelf dataset [20], which is crawled from Twitter through its open API. The details about the dataset can be summarized as Table 1.

**Table 1.** Twitter Dataset Statistics

| | |
|---|---|
| Total users 139,180 | Friends per user 14.8 |
| Total edges 4,175,405 | Followers per user 14.9 |
| Total tweets 76,409,820 | Tweets per user 549 |

It is time-consuming to establish subjectivity model with the 139,180 users directly for the computational complexity of LDA. However, the principle of homophily [21], or "birds of a feather flock together" [22] suggests that users that are "connected" closely may tend to talk about similar topics and hold similar opinions [23]. On Twitter, the connections a user creates may correspond to approval or a desire to pay

attention, or suggestive of the possibility of common topics and opinions. Therefore we adopt the community structure of social network to divide the 139,180 users into different community and establish subjectivity model for a user in his community local network. The communities are found with the packages igraph[1]. There are 106 communities in the global network, and 73 communities consist of users less than 15, for which topics can't be found effectively with LDA, so we filter out users in these communities. At the same time, we also filter out 15,756 users who are inactive with tweets less than 5, only tweet themselves with words less than 3, or only publish content with url links. In the final dataset, there are 122,329 users distributed in 33 communities. The subjectivity model for each user is established within his own community as algorithm 1.

Besides our model, we also conduct a set of experiments comparing with other topic-sentiment model including JST and TSM. The symmetry Dirichlet priors of topic models were set to 50/T and 0.01 respectively. The asymmetry sentiment prior empirically was set to (0.01, 1.8) for JST. All results were averaged over 5 runs with 2000 Gibbs sampling iterations.

## 5.2   Case Study

In order to qualitatively evaluate the effectiveness of our method, we give a vivid example of a user's subjectivity model, who has published 533 tweets. All his tweets are illustrated as Figure 3(a) in a word-cloud figure.
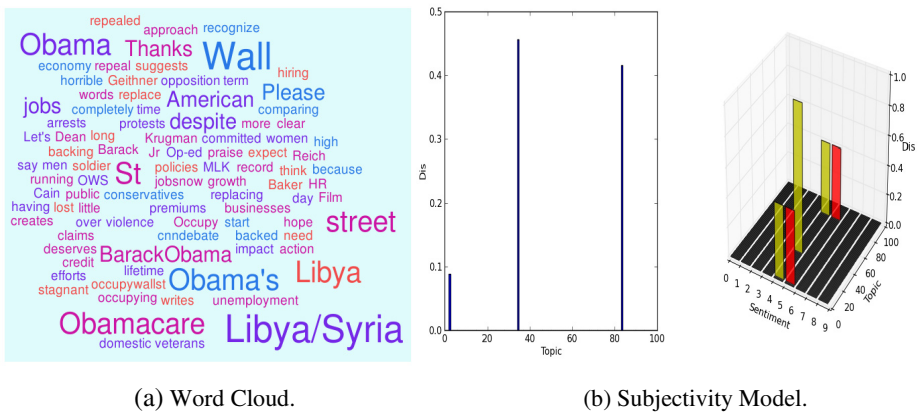


(a) Word Cloud.                      (b) Subjectivity Model.

**Fig. 3.** Example of an user. In the subjectivity model, left sub-graph 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 sub-graph 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).

Figure 3(b) is the visualized subjectivity model of the user in a [0, 100] topic space and a [0, 8] sentiment space, which is established according to our method. It is obvious that the user is interested in three topics (topic 2: "#Obamacare", topic 32:

---

[1]   http://igraph.org/

"#libya" and topic 83:"#occupywallst"), and the left part of Figure 3(b) denotes the weights of his topics of interest. The right part denotes the opinions of the user towards three topics, in which he is neutral to topic "#libya" with 100% distribution on sentiment strength value 4, positive to topic "#Obamacare" and "#occupywallst" with 50% on value 4 and 50% on value 5. From the example, it is demonstrated that our model can give a detail description for the subjectivity of users in that it can model not only the interest distribution but also opinion coverage over a fine-grained sentiment.

## 5.3    Opinion Prediction Performance

To directly evaluate the effectiveness of our model quantitatively, we compare our model with other two generative topic-sentiment model (TSM and JST) with the number of topic is set to 50, 100, 150 and 200 iteratively. Short of labeled training data, we only compare our method with three state-of-the-art unsupervised sentiment analysis methods in the performance of opinion prediction.

- OF: OpinionFinder is a publicly available software package for sentiment analysis that can be applied to determine sentence-level subjectivity, i.e. to identify the emotional polarity (positive or negative) of sentences [24].
- S140: Sentiment140 can automatically classifying the sentiment of tweets using distant supervision with training data consisted of Twitter messages with emoticons.
- STR: SentiStrength package has been built especially to cope with sentiment analysis in short informal text of social media. It combines lexicon-based approaches with sophisticated linguistic rules adapted to social media [25].

We randomly select 1,000 target users from our dataset with at least 80 tweets, and select one random tweet for each user from his tweets collection to form a set of 1,000 tweets for evaluation. In order to identify topic of each tweet easily, the tweets with hashtag are prior to be selected. All 1,000 tweets in the test set are manually labeled with sentiment polarity as the golden standard. Accuracy is used as our performance measurement, and the result is list in Table 2.

**Table 2.** Accuracy performance. A significant improvement over OF with*

| Method | 50 | 100 | 150 | 200 |
|--------|------|------|------|------|
| OF | 65.85% | | | |
| S140 | 70.45%* | | | |
| STR | 69.98%* | | | |
| TSM | 63.46% | 72.94%* | 67.83% | 66.65% |
| JST | 61.25% | 68.57%* | 75.88%* | 67.03% |
| SUB | 71.53%* | 81.05%* | 78.32%* | 74.54%* |

As can be observed from the result table that:

Firstly, the performance of OpinioFinder is the lowest with 65.85% accuracy, and we think the reason lying in that it is designed for the review and not adapts to tweets with informal language usage;

Secondly, other two unsupervised sentiment methods (Sentiment140: 70.45% and SentiStrength: 69.98%) outperform OpinioFinder significantly.

Thirdly, overall, the two generative models outperform OpinionFinder significantly, which demonstrates the importance of relating sentiment to the topics of users. Their performances are a little better than Sentiment140 and SentiStrength, but not significantly.

Finally, our method (SUB) outperforms all three unsupervised sentiment methods significantly with all four topic settings, and improves the performance of Senti-Strength significantly by combining subjectivity model of users with content of tweet. Compared with two generative models, our model outperforms TSM significantly, and gets a little better performance than JST. We think it is because sentiment analysis technique of our model is more suitable for the Twitter language, for it can extract subtle sentiment imbedded in special language characteristics such as repeated letters and emoticons.

## 6    Conclusion

In this paper, we define and investigate a novel opinion integration problem for social media users. We propose a subjectivity model to solve this problem in a three-stage framework and design an algorithm to establish the subjectivity model from historical tweets of users. With this model, we can automatically generate an integrated opinion summary that consists of both topics of interest distribution and topic related opinion distribution for a user. The proposed model is demonstrated effective in the application of opinion prediction. Experiments on Twitter data show that the proposed model can effectively describe topic related opinions with two probabilistic distributions and clearly outperforms generative models in the opinion prediction task. In the future, we will apply our model in several social network analysis applications to testify its effectiveness.

## References

1. Lu, Y., Zhai, C.: Opinion integration through semi-supervised topic modeling. In: Proceedings of the 17th International Conference on World Wide Web, pp. 121–130. ACM (2008)
2. Pang, B., Lee, L.: Opinion mining and sentiment analysis. Foundations and Trends in Information Retrieval 2(1-2), 1–135 (2008)
3. Liu, B.: Sentiment analysis and opinion mining. Synthesis Lectures on Human Language Technologies 5(1), 1–167 (2012)
4. Barbosa, L., Feng, J.: Robust sentiment detection on twitter from biased and noisy data. In: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, pp. 36–44. Association for Computational Linguistics (2010)
5. Davidov, D., Tsur, O., Rappoport, A.: Enhanced sentiment learning using twitter hashtags and smileys. In: Proceedings of the 23rd International Conference on Computational Linguistics: Posters, pp. 241–249. Association for Computational Linguistics (2010)
6. Jiang, L., Yu, M., Zhou, M., Liu, X., Zhao, T.: Target-dependent twitter sentiment classification. In: Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, vol. 1, pp. 151–160. Association for Computational Linguistics (2011)

7. Li, G., Hoi, S.C., Chang, K., Jain, R.: Micro-blogging sentiment detection by collaborative online learning. In: 2010 IEEE 10th International Conference on Data Mining (ICDM), pp. 893–898. IEEE (2010)

 8. Tan, C., Lee, L., Tang, J., Jiang, L., Zhou, M., Li, P.: User-level sentiment analysis incorporating social networks. In: Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1397–1405. ACM (2011)

 9. Mostafa, M.M.: More than words: Social networks text mining for consumer brand sentiments. Expert Systems with Applications 40(10), 4241–4251 (2013)

10. Malouf, R., Mullen, T.: Taking sides: User classification for informal online political discourse. Internet Research 18(2), 177–190 (2008)

11. Liu, H., Zhao, Y., Qin, B., Liu, T.: Comment target extraction and sentiment classification. Journal of Chinese Information Processing 24(1), 84–89 (2010)

12. Zhai, Z., Liu, B., Xu, H., Jia, P.: Constrained LDA for grouping product features in opinion mining. In: Huang, J.Z., Cao, L., Srivastava, J. (eds.) PAKDD 2011, Part I. LNCS (LNAI), vol. 6634, pp. 448–459. Springer, Heidelberg (2011)

13. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. The Journal of Machine Learning Research 3, 993–1022 (2003)

14. Rosen-Zvi, M., Griffiths, T., Steyvers, M., Smyth, P.: The author-topic model for authors and documents. In: Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence, pp. 487–494. AUAI Press (2004)

15. Ramage, D., Hall, D., Nallapati, R., Manning, C.D.: Labeled lda: A supervised topic model for credit attribution in multi-labeled corpora. In: Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing, vol. 1, pp. 248–256. Association for Computational Linguistics (2009)

16. Mei, Q., Ling, X., Wondra, M., Su, H., Zhai, C.: Topic sentiment mixture: modeling facets and opinions in weblogs. In: Proceedings of the 16th International Conference on World Wide Web, pp. 171–180. ACM (2007)

17. Lin, C., He, Y.: Joint sentiment/topic model for sentiment analysis. In: Proceedings of the 18th ACM Conference on Information and Knowledge Management, pp. 375–384. ACM (2009)

18. Xie, S., Tang, J., Wang, T.: Resonance elicits diffusion: Modeling subjectivity for retweeting behavior analysis. Cognitive Computation, 1–13 (2014)

19. Walton, D.N.: Bias, critical doubt and fallacies. Argumentation and Advocacy 28, 1–22 (1991)

20. Li, R., Wang, S., Deng, H., Wang, R., Chang, K.C.C.: Towards social user profiling: unified and discriminative influence model for inferring home locations. In: KDD, pp. 1023–1031 (2012)

21. Lazarsfeld, P.F., Merton, R.K.: Friendship as a social process: A substantive and methodological analysis. In: Berger, M., Abel, T. (eds.) Freedom and Control in Modern Society. Van Nostrand, New York (1954)

22. McPherson, M., Smith-Lovin, L., Cook, J.M.: Birds of a feather: Homophily in social networks. Annual Review of Sociology, 415–444 (2001)

23. Thelwall, M.: Emotion homophily in social network site messages. First Monday 15(4) (2010)

24. Wilson, T., Wiebe, J., Hoffmann, P.: Recognizing contextual polarity in phrase-level sentiment analysis. In: Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing, pp. 347–354. Association for Computational Linguistics (2005)

25. Thelwall, M., Buckley, K., Paltoglou, G., Cai, D., Kappas, A.: Sentiment strength detection in short informal text. Journal of the American Society for Information Science and Technology 61(12), 2544–2558 (2010)