# Confluence: Conformity Influence in Large Social Networks

Jie Tang
Computer Science
Tsinghua University, China
Beijing 100084, China
jietang@tsinghua.edu.cn

Sen Wu
Computer Science
Tsinghua University, China
Beijing 100084, China
ronaldosen@gmail.com

Jimeng Sun
IBM T. J. Watson Research
Center
USA
jimeng@us.ibm.com

## ABSTRACT

Conformity is a type of social influence involving a change in opinion or behavior in order to fit in with a group. Employing several social networks as the source for our experimental data, we study how the effect of conformity plays a role in changing users' online behavior. We formally define several major types of conformity in individual, peer, and group levels. We propose *Confluence* model to formalize the effects of social conformity into a probabilistic model. Confluence can distinguish and quantify the effects of the different types of conformities. To scale up to large social networks, we propose a distributed learning method that can construct the Confluence model efficiently with near-linear speedup.

Our experimental results on four different types of large social networks, i.e., Flickr, Gowalla, Weibo and Co-Author, verify the existence of the conformity phenomena. Leveraging the conformity information, Confluence can accurately predict actions of users. Our experiments show that Confluence significantly improves the prediction accuracy by up to 5-10% compared with several alternative methods.

## Categories and Subject Descriptors

J.4 [Social and Behavioral Sciences]: Miscellaneous; H.3.3 [Information Search and Retrieval]: Text Mining

## General Terms

Algorithms, Experimentation

## Keywords

Conformity; Social influence; Social network

## 1. INTRODUCTION

Conformity is the act of matching attitudes, beliefs, and behaviors to group norms [8]. The phenomenon of conformity could occur in small groups or the whole society, as a resultant of peer influence or group pressure. Conformity can have either good or bad effect depending on the situation. For example, it helps form and maintain the social norms, and helps prevent acts that are perceptually dangerous. Conformity can be influenced by various factors such as individual status, peer influence and group pressure. Therefore there is a clear need for quantitative methods for measuring conformity from different aspects, so as to understand the complex dynamics in social networks.
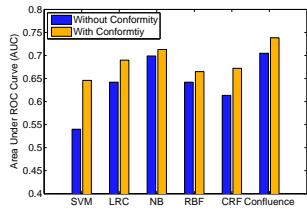
Conformity was first studied by psychologists through interviews with small groups of participants [17]. In economics, Bernheim [3] found that sometimes people are willing to conform simply because they recognize that departure from the social norm may impair their status. Bernheim further proposed a theory of social conformity and presented a model to describe the conformity process. However, due to the lack of real data, he only studied the model from the theoretical aspect. With the rapid proliferation of online social networks such as Facebook, Twitter, and Flickr, it becomes feasible and also very necessary to conduct an in-depth investigation of the conformity problem on real large social networks. In practice, the effect of conformity has been also observed in online social networks. For example, Bond et al. [4] reported results from a randomized controlled trial of political mobilization messages delivered to 61 million Facebook users. They found that when one is aware that their friends have made the political votes, their likelihood to vote will significantly increase. Bakshy et al. [2] also found that when their friends click an ad, they will be more likely to click the same ad.

From a broader viewpoint, conformity can be seen as a special type of social influence. There are a bulk of studies on social influence analysis. These studies can be roughly classified int
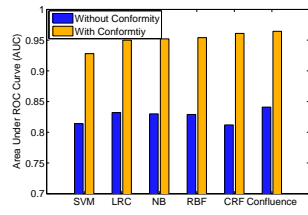
(a) Flickr    (b) Gowalla



SVM

*De nition* **Group on or _ t** The group conformity is then defined as the ratio between the number of actions for which we have evidence that the user $v$ conforms to the group, over the total number of $\tau$-group actions performed by users in the group $C_k$,

$$gcf^\tau(v, C_{vk}) = \frac{|(a, v', t') \in A^\tau_{C_k} | \exists (a, v, t) : \mathbb{I}[c_{ik}] \wedge \epsilon \geq t - t' \geq 0|}{|A^\tau_{C_k}|}$$

where $A^\tau_{C_k} \subset A$ denotes actions performed by more than a percentage $\tau$ of all users in the group $C_k$; $\mathbb{I}[c_{ik}]$ is an indicator function, which returns true if the value of $c_{ik}$ is 1 and false otherwise.

Please note that a user may be involved into more than one groups, thus has different conformity degrees in the different groups. The above definitions quantify the conformity from differ-

normalization factor to ensure that the distribution is normalized so that the sum of the probabilities equals to 1.

In practice, choosing a good threshold $\epsilon$ for defining the conformity factors is challenging. Instead, we use a decay factor $\lambda \geq 1$ in each conformity function. A large $\lambda$ means a slow-decay effect. Accordingly, the peer conformity factor is defined as:

$$g(y_i, y'_j, pcf(v_i, v_j)) = (\frac{1}{2})^{\frac{t-t'}{\lambda}} pcf(v_i, v_j) \qquad (2)$$

where $t'$ corresponds to the latest past time when $v_j$ performed the same action as $v_i$ in the training data set. The decay factor decays the peer conformity exponentially over time, with the half-life, $\lambda$, serving as a tunable parameter. The basi0

We first introduce how we calculate the gradient for each parameter. As the network structure in the social network can be arbitrary (may contain cycles), it is intractable to obtain exact solution of the objective function using methods such as Junction Tree [34]. We use Loopy Belief Propagation (LBP) [26] to approximate the solution. Specifically, we first approximate the marginal distribution $P_\theta(y_i|.)$ using LBP. With the marginal probabilities, the gradient can be obtained by summing over all factor functions. Theoretically, the LBP algorithm does not guarantee a convergence and may result in local maximum, but in practice its performance is good. We empirically compare the effectiveness and efficiency of the algorithm in Section 4. After obtained the marginal distribution $P_\theta(y_i|.)$, we use a gradient descent method (or a Newton-Raphson method) to solve the objective function (Eq. 1). We use $\alpha$ as the example to explain how we learn the parameters. Specifically, we first write the gradient of each unknown parameter $\alpha$ with regard to the objective function:

$$\frac{\mathcal{O}(\theta)}{\alpha_j} = \mathbb{E}[f(y_i, x_{ij})] - \mathbb{E}_{P(y_i|G,A)}[f(y_i, x_{ij})] \quad (5)$$

where $\mathbb{E}[f(y_i, x_{ij})]$ is the expectation of the local factor function $f(y_i, x_{ij})$ given the data distribution in the input network and $\mathbb{E}_{P(y_i|G,A)}[g(y_i, x_{ij})]$ represents the expectation under the distribution learned by the model, i.e., $P(y_i|G, A)$. Similar gradients can be derived for parameter $\beta_i$, $\gamma_{ij}$, and $\mu_{ik}$.

Now we explain how we use distributed learning to approximate the marginal probability. We use a master-slave architecture, i.e., one master machine is responsible for optimizing parameters, and the other slave machines are responsible for calculating the marginal probabilities. At the beginning of the algorithm, the graphical model of Confluence is partitioned into $M$ roughly equal subgraphs, where $M$ is the number of slave processors. The partition can be done by any graph cut software. After the partition, the subgraphs are then distributed over slave processors. Then each slave processor calculates the "local" belief (the marginal probability) on the subgraph $G_l$ according to the following equations (again we use $P(y_i|G, A)$ as the example in the explanation):

$$m_{ij}^l(y_i) = \sigma \sum_{y_i} \psi_{ij}^l(y_i, y_j) \psi_i^l(y_i) \prod_{k \in NB(i) \setminus j} m_{ki}^l(y_i) \quad (6)$$

$$b_i^l(y_i) = \psi_i^l(y_i) \prod_{k \in NB(i)} m_{ki}^l(y_i) \quad (7)$$

$$P(y_i|.) = \sigma \sum_{l=1}^{M} b_i^l(y_i) \quad (8)$$

where $\sigma$ denotes a normalization constant; $m_{ij}^l(y_i)$ is the "belief" propagated from node $y_j$ to node $y_i$; $NB(i) \setminus j$ denotes all nodes neighboring node $y_i$ in the subgraph $G_l$, except $y_j$; $\psi_i^l(y_i)$ denotes all defined factor functions related to $y_i$ in the subgraph $G_l$ and is calculated by $\psi_i^l(y_i) = \exp(\sum_{k=1}^{d} f(y_i, x_{ik}) + \beta_i g(y_i, icf(v_i)))$, and $\psi_i^l(y_i, y_j)$ denotes all correlation factor functions related to $y_i$ in the subgraph; notation $b_i^l(y_i)$ denotes the unnormalized "local" belief collected from each subgraph, and finally by combining them together we obtain the marginal probability $P(y_i|.)$.

However, inevitably there will be some correlation factors defined over nodes that are partitioned into different subgraphs. These correlation factors cannot be calculated due to the high communication cost. Simply eliminating those correlation factors may hurt the learning performance. To alleviate this problem, we present a virtual node based method. In particular, suppose three nodes $(y_1, y_2, y_3)$ in the Confluence model are associated with a group

conformity factor $g(.)$. If the partition assigns two nodes (e.g., $y_1$ and $y_2$) into one subgraph $G_1$ and the rest one (i.e., $y_3$) into another subgraph $G_2$, then we create a virtual node in the first subgraph $G_1$ so that the group conformity factor can be still calculated in the subgraph. For the virtual node, we do not calculate the local attribute factors $f(.)$. The distributed learning algorithm is summarized in Algorithm 1.

**o    n r n**    The learned model parameters $\theta$ can be used to infer users' future actions. In particular, given the network $G$ and the action history $A$, we aim to predict users action labels $Y^{t+1}$ at time $t + 1$. This can be done by performing the model inference on the network to maximize the conditional probability, i.e.,

$$Y^\star = \arg\max_{Y^{t+1}} P_\theta(Y^{t+1}|G, A) \quad (9)$$

Again, we use the distributed loopy belief propagation algorithm to compute the marginal probability $P_\theta(y_i^{t+1}|.)$ and then predict the action of each user at time $t + 1$ as the label that has the largest marginal probability. For each user, we define the individual conformity factor according to the estimated individual conformity from the training data. For defining the peer conformity factor $g(y_i, y_j', pcf(v_i, v_j))$ between $v_i$ and $v_j$, we first find the latest past time $t'$ when $v_j$ performed the corresponding action $y_j'$, and calculate the factor according to Eq. 2. The group conformity factor can be similarly defined.

## 4  EX  L E A  ES

We conduct various experiments to evaluate the Confluence method. The datasets and codes are publicly available.[3]

## 4    E p r  nt S tup

**D t  s ts**    We evaluate the proposed method on four different genres of networks: Flickr, Gowalla, Weibo, and Co-Author. Table 2 lists statistics of the four networks.

**F    r** is a photo sharing network. Users on the site can share photos and add comments to other photos. Flickr users can also create and join different groups. The data set spans the period from Apr. 1st, 2012 to Jun. 16th, 2012. We define the action as adding a comment to a specific photo. Thus the action space includes all photos on Flickr. To avoid the sparsity problem, we remove those photos with less than 5 comments. This results in 144,627 unique actions. We try to study how users' commenting actions conform to the other users in the network.

**Go** is a location-based social network, where users share their locations by checking-in. The data was from [7] and all check-ins of these users over the period of Jul. 10th, 2010 - Jul. 29th, 2010. The action in this data set is defined as whether a user checks in some location (indicated by hashtag or location ID). Thus the dimension of the action space is the number of available locations. We also remove those locations with less than five check-ins and finally obtain 218,811 unique actions. Our goal is to study whether the users will conform to their friends' check-in behavior.

**o** is the most popular microblogging service in China. We collected a complete network between 1,700,000 users and all the tweets posted by those users between Sep. 28th, 2012 and Oct. 29th, 2012. The action is defined as whether a user posts a message on a specific topic (indicated by hashtag). We choose the ten most popular topics in 2012 and study how users conform to each other in the network on discussing those topics. We aim to study how users conform to each other on discussing those topics.

---

[3]http://arnetminer.org/conformity/

Table 2: Statistics of the four networks

| Dataset | Flickr | Gowalla | Weibo | Co-Author |
|---|---|---|---|---|
| #nodes | 1,991,509 | 196,591 | 1,776,950 | 737,690 |
| #edges | 208,118,719 | 950,327 | 308,489,739 | 2,416,472 |
| #groups | 460,888 | N/A | N/A | 60 |
| #actions | 3,531,801 | 6,442,890 | 6,761,186 | 1,974,466 |



Figure 4: Factor Contribution Analysis. Confluence_base stands for training Confluence_base (without the conformity structures). Confluence_base+SB stands for Confluence_base with social-balance structures; Confluence_base+SH stands for Confluence_base with structural-hole structures; Confluence_base+OL stands for Confluence_base with opinion-leader structures; Confluence_base+ST stands for Confluence_base with strong-tie structures; Confluence_base+CF stands for Confluence_base with conformity structures.

**Co Author** is a network of authors. The data set, crawled by Arnetminer.org [33], is comprised of 737,690 CS authors and 2,416,472 co-authorships over 1975 - 2012. Based on the publication venues, authors are categorized into different domains such as Data Mining, Artificial Intelligence, Computer Graphics, etc.[4] The action is defined as whether an author will publish a paper in a specific domain. Thus, in total we have 200 unique actions. Our goal is to study how an author conforms to the other authors on choosing the publication venue.

**Evaluation Metrics** To quantitatively evaluate the proposed model, we consider the following performance metrics:

- **Action prediction** We apply the learned model for action prediction and evaluate its performance in terms of Precision, Recall, F1-Measure, and Area Under Curve (AUC).

- **Scalability performance** We evaluate the computational time as the efficiency metric.

- **Qualitative case study** We use a case study as the anecdotal evidence to further demonstrate the effectiveness of the proposed model.

All codes are implemented in C++ and JAVA, and all the evaluations are performed on an x64 machine with E7520 1.87GHz Intel Xeon CPU (with 16 cores) and 192GB RAM. The operation system is Microsoft Windows Server 2008 R2 Enterprise. The proposed distributed learning algorithm has a good convergence property. On average, it converges within 100 iterations.

**Comparison Methods** Given the input network $G$ and the action history $A$, we can construct a training data set $\{(x_i, y_i)\}_{i=1,\cdots,n}$, where $n = |A|$; $x_i$ is the feature vector associated with user $v_i$ and $y_i = a$ indicates whether user $v_i$ performs the corresponding action $a$. In this way, we can use existing methods such as Support Vector Machines (SVMs) or Logistic Regression (LR) to train a classification model and then apply the trained model to predict users' future actions. The difference from our proposed Confluence model is that the classification model does not consider the correlation between users' actions. We also compare with Conditional Random Fields (CRFs) [21].

**SVM** it uses all defined features associated with each user to train a classification model and then apply it to predict users' actions in the test data. For SVM, we use SVM-light.[5]
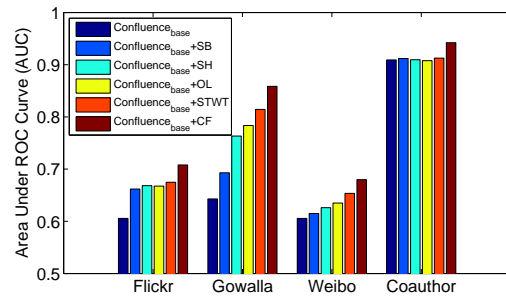
**LR** it uses logistic regression (LR) to train the classification model with the same features as those in the SVM method. We also compare with the results of Naive Bayes (NB) and Gaussian Radial Basis Function Neural Network (RBF). For all the three methods, we employ Weka.[6]

**CRF** it is a graphical model based on Conditional Random Field (CRF). Comparing with CRFs, the factor graph model provides a

---

[4]Refer to http://arnetminer.org/topic-browser for a list of domains.
[5]http://svmlight.joachims.org/
[6]http://www.cs.waikato.ac.nz/ml/weka/

more explicit explanation for the factorization of the underlying probability distribution [19]. In addition, it is not easy to incorporate the group conformity factors into the CRF model, as users' group memberships could be arbitrary and one user can belong to multiple groups. Thus in the CRF method, we use attribute-based features, the social-based features, and individual conformity features, but do not use the group conformity features. For CRF, we use Mallet [25].

In all the comparison methods, we try to use the same features. The attribute based features are used in all the methods. The social features defined on social ties and social balance are used in CRF and Confluence only, as SVM and LR cannot capture the correlations. As for the conformity features, individual conformity is defined for each user, and is used in all methods; peer conformity is defined for peer friends and is used in CRF and Confluence; group conformity is defined for groups and is used only in Confluence.

## 4 Action Prediction Analysis

On all the four data sets, we use the historic users' actions as the training data in different methods and use the learned model to predict users' action in the next time stamp. Specifically on Flickr each week is a time stamp (which results in 11 time stamps in total), on Gowalla and Weibo each day a time stamp (which result in 20 and 32 time stamps respectively), and on Co-Author each year is a time stamp (which results in 38 time stamps). We perform the prediction for each time stamp and finally report the average performance.

**Action prediction performance** Table 3 lists the action prediction performance of the different methods on the four data sets. Our method Confluence consistently achieves better performance than the comparison methods. In terms of F1-score, Confluence achieves a 1-17% improvement compared with the SVM, LR, NB, and RBF methods that do not consider the correlation features. CRF also considers some correlation features (such as social tie and social balance based features), thus improves the prediction performance. However it cannot incorporate the group conformity feature, thus still underperforms our method by 0.5-10.5% in terms of F1-score. We produced sign tests for each result, which confirms that all the improvements of our proposed models over the five methods are statistically significant ($p \ll 0.01$).

**Factor contribution analysis** In the Confluence model, we define basic features based on the user-associated attributes, and five

| Data | Method | Precision | Recall | F1-measure | AUC |
|------|--------|-----------|--------|------------|-----|
| **Flickr** | SVM | 0.5921 (±0.0036) | 0.5905 (±0.0031) | 0.5802 (±0.0012) | 0.6473 (±0.0004) |
| | LR | 0.6010 (±0.0052) | 0.5900 (±0.0057) | 0.5770 (±0.0018) | 0.6510 (±0.0008) |
| | NB | 0.6170 (±0.0071) | 0.6040 (±0.0083) | 0.5920 (±0.0031) | 0.6520 (±0.0019) |
| | RBF | 6 5 ± | 0.5960 (±0.0010) | 0.5720 (±0.0024) | 0.6700 (±0.0010) |
| | CRF | 0.5474 (±0.0030) | ± | 0.6239 (±0.0016) | 0.6722 (±0.0010) |
| | Confluence | 0.5472 (±0.0025) | 0.7770(±0.0010) | 6 4 ± | ' ± 6 |
| **Gowalla** | SVM | 0.9290 (±0.0212) | 0.9310 (±0.0121) | 0.9295 (±0.0105) | 0.9280 (±0.0042) |
| | LR | 0.9320 (±0.0234) | 0.9290 (±0.0234) | 0.9310 (±0.0155) | 0.9500 (±0.0054) |
| | NB | 0.9310 (±0.0197) | 0.9290 (±0.0335) | 0.9300 (±0.0223) | 0.9520 (±0.0030) |
| | RBF | 0.9320 (±0.0254) | 0.9280 (±0.0284) | 0.9300 (±0.0182) | 0.9540 (±0.0022) |
| | CRF | 0.9330 (±0.0100) | 0.9320 (±0.0291) | 0.9330 (±0.0164) | 0.9610 (±0.0019) |
| | Confluence | ' ± ' | ± ' | 5 ± | 644 ± 4 |
| **Weibo** | SVM | 0.5060 (±0.0381) | 0.5060 (±0.0181) | 0.5060 (±0.0157) | 0.5070 (±0.0053) |
| | LR | 0.5190 (±0.0461) | 0.6450 (±0.0104) | 0.5750 (±0.0281) | 0.5390 (±0.0133) |
| | NB | 0.5120 (±0.0296) | 0.6700 (±0.0085) | 0.5810 (±0.0165) | 0.5390 (±0.0132) |
| | RBF | 5 4 ± 4 | 0.5690 (±0.0098) | 0.5460 (±0.0159) | 0.5450 (±0.0103) |
| | CRF | 0.5150 (±0.0353) | 0.6310 (±0.0121) | 0.5720 (±0.0209) | 0.6320 (±0.0139) |
| | Confluence | 0.5185 (±0.0296) | 6 ± 5 | 6 6 ± 56 | ' 5 ± '' |
| **Co-Author** | SVM | 0.7672 (±0.0338) | 0.8671 (±0.0145) | 0.8256 (±0.0129) | 0.8562 (±0.0115) |
| | LR | 0.8700 (±0.0261) | 0.7640 (±0.0346) | 0.8140 (±0.0221) | 0.8500 (±0.0030) |
| | NB | 0.7640 (±0.0177) | 0.8510 (±0.0185) | 0.8050 (±0.0048) | 0.8720 (±0.0074) |
| | RBF | 0.7720 (±0.0182) | 0.8830 (±0.0191) | 0.8240 (±0.0145) | 0.8790 (±0.0031) |
| | CRF | 0.8081 (±0.0252) | 0.8771 (±0.0249) | 0.8360 (±0.0087) | 0.9025 (±0.0025) |
| | Confluence | ± 5 | ± | ± 4 | 5 ± |

types of social features: social balance (SB), structure hole (SH), opinion leader (OL), strong tie/weak tie (STWT), and conformity (CF). Here we examine the contributions of the different social factors defined in our model. Specifically, first we use the basic features to train a model (referred to as Confluence$_{base}$). Then we incrementally add one of the five social features and evaluate its improvement on the prediction performance over that using only basic features. Figure 3 shows the Area Under Curve (AUC) score on the different data sets. We see that different social factors contribute differently in the different networks. For example, the opinion leader based features are very useful in the Gowalla network, but less useful in the Co-Author network. On the other hand, the conformity based features consistently improve the prediction performance on all the networks. In terms of the AUC score, the improvements by adding conformity features range from 2% to 20% in the four networks. This analysis confirms the importance of the conformity phenomena in social networks.

**Effects of conformity** We further present an in-depth analysis of how different levels of conformities affect the performance of action prediction. Figure 4 shows the prediction performance (in terms of AUC) of the proposed Confluence by considering different levels of conformities. Confluence$_{base}$ stands for the Confluence method by considering only basic features (i.e., ignoring all conformity factors). It can be clearly seen that without the conformity based factors, the prediction performance drop significantly. Co-Author network is most predictable because the co-authorships are stable and predictable in general. Weibo and Flickr are the most difficult to predict because the user behavior is fairly autonomous and independent. Conformity has most significant prediction impact on Gowalla, which suggests conformity plays an important role in geospatial and mobile applications in social networks. By incorporating the conformity features, significant improvements (+20-30%) over the prediction performance can be ob-
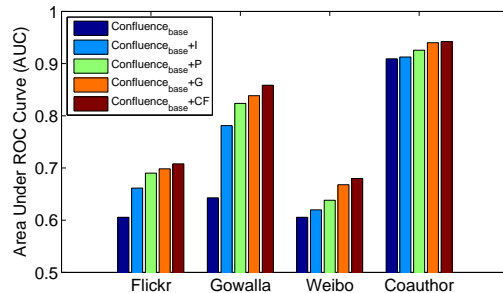


Figure 4. Effects of conformity. Confluence$_{base}$ stands for the Confluence method without using social conformity features. Confluence$_{base}$+I stands for the Confluence$_{base}$ method plus using individual conformity features. Confluence$_{base}$+P stands for the Confluence$_{base}$ method plus using peer conformity features. Confluence$_{base}$+G stands for the Confluence$_{base}$ method plus using group conformity features.

tained on Gowalla. Confluence$_{base}$+I (or +P or +G) respectively indicates that we respectively add individual conformity features (or peer conformity features or group conformity features) into the Confluence$_{base}$ method. By incorporating each type of conformity factors, we observe clear improvement compared to the Confluence$_{base}$ method. We can also see that on all the four data sets, the group conformity is more important than the other two types of conformities. This makes sense, as in most cases conformity is a group phenomenon rather than an individual behavior.

## 4. Scalability Performance

We now evaluate the scalability performance of the distributed learning algorithm on the four networks. In our experiments, we
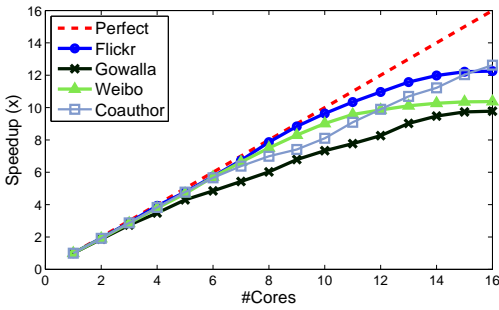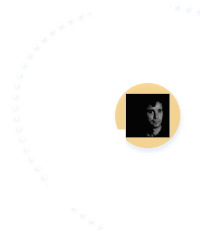
Figure 5: Speedup of distributed learning

Table 4: Running time of the proposed method (hour)

| Data Set | Flickr | Gowalla | Weibo | Co-Author |
|---|---|---|---|---|
| Confluence | 1.602 | 0.245 | 1.083 | 0.512 |
| Confluence (single) | 19.637 | 2.395 | 11.229 | 6.464 |
| CRF | 3.864 | 0.387 | 2.547 | 1.823 |

use METIS [16] to partition the graph into multiple subgraphs (one for each core). Figure 5 shows the speedup of the distributed algorithm with different number of computer nodes (2, 3, 6, 8, 10, 12, 14, 16 cores) used. The speedup curve is close to the perfect line at the beginning. Though the speedup inevitably decreases due to the increase of the communication cost between the different computer nodes, the distributed learning algorithm can still achieve $\sim 9\times$ speedup with 16 cores. It is noticeable that the speedup curves on different networks present a bit different patterns. This is due to the difference of the network properties (such as densities). Table 4 further gives the running time for learning the proposed Confluence model over 16 computer cores and single compute on different data sets.

Another thing worth noting is that the distributed learning is essentially an approximation of the original learning algorithm on a single machine. We used METIS to partition the graph into multiple subgraphs and distribute the subgraphs onto slave machines. We also evaluate the prediction performance by the distributed learning algorithm. On average, the prediction accuracy by the distributed learning over 16 cores only drops slightly (ranging from 0.5-1.68%), which further demonstrates the effectiveness of the distributed learning algorithm.

## 4.4 About the Case Study

Now we use a case study from Flickr to further demonstrate the effectiveness of the proposed model. Figure 6 shows an example extracted from Flickr. User A joined three groups (denoted as Group 1, 2, 3 respectively). On 03/10/2012, user A added one comment respectively to Picture 1 and Picture 2. The Action 1 (adding comment to Picture 1) was mainly performed in Group 1 and the Action 2 (adding comment to Picture 2) was mainly performed in Group 2. After modeling with the proposed Confluence method, the modeling results suggest that, for performing Action 1, user A has a strong conformity to user B, but very weak conformity to user D and C. By taking a closer look at the data, we found that Group 1 is a loosely connected group and members have very few connections in the group, and the comments to the same photo are very controversial (such as the comments of B and D to Picture 1). Thus the influence between users are mainly at the peer level. For Ac-

to study how user behavior is influenced by close friends in their ego networks. Li et al. [22] tried to study the interplay between influence and individual conformity. However, they do not consider the group conformity. Quite a few studies have been done for maximizing the influence spread in social network. Domingos and Richardson [10, 28] formally defined influence maximization as an algorithmic problem and prove its NP-hardness. Chen et al. [6] further developed efficient algorithms to approximately solve the influence maximization problem. While, influence maximization is in nature different from the conformity analysis problem. To the best of our knowledge, this is the first attempt to formally define the problem of conformity influence analysis and to address this problem with a principled method.

## 6 CONCLUSION

In this paper, we study a novel problem of conformity influence analysis in large social networks. We formally define three major types of conformities, precisely formulate the problem of conformity influence analysis, and propose a Confluence model to model users' actions and conformity. Three factor functions are defined to capture the different levels of conformities. A distributed learning algorithm is presented to efficiently learn the proposed model. We validate the effectiveness and efficiency of the proposed model on four networks. Our experimental results show that the proposed method significantly outperforms several alternative methods. We also present a case study to further demonstrate the effectiveness of the method.

Understanding the fundamental mechanism of social conformity is very important for social network analysis and represents a new and interesting research direction. As for the future work, it would be intriguing to connect the conformity phenomenon with some other social theories such as social status and structural holes so as to understand the formation and dynamic change of the network structure. It is also interesting to design some other model, for example a game theory based model, to model the conformity phenomenon. As for the proposed Confluence model itself, it has many parameters. We also consider adding regularization to control the sparsity of those parameters.

## REFERENCES

[1] A. Anagnostopoulos, R. Kumar, and M. Mahdian. Influence and correlation in social networks. In *KDD 0*, pages 7–15, 2008.

[2] E. Bakshy, D. Eckles, R. Yan, and I. Rosenn. Social influence in social advertising: evidence from field experiments. In *EC*, pages 146–161, 2012.

[3] B. D. Bernheim. A theory of conformity. *Journa of Po itica Econo y*, 1027(5):841–877, 1994.

[4] R. M. Bond, C. J. Fariss, J. J. Jones, A. D. I. Kramer, C. Marlow, J. E. Settle, and J. H. Fowler. A 61-million-person experiment in social influence and political mobilization. *Nature*, 489:295–298, 2012.

[5] R. S. Burt. *Structura Ho es The Socia Structure of Co petition*. Harvard University Press, 1992.

[6] W. Chen, Y. Wang, and S. Yang. Efficient influence maximization in social networks. In *KDD 0*, pages 199–207, 2009.

[7] E. Cho, S. A. Myers, and J. Leskovec. Friendship and mobility: user movement in location-based social networks. In *KDD*, pages 1082–1090, 2011.

[9] D. Crandall, D. Cosley, D. Huttenlocher, J. Kleinberg, and S. Suri. Feedback effects betw