# Happiness Is Assortative in Online Social Networks

Johan Bollen*,**
Indiana University

Bruno Gonçalves**
Indiana University

Guangchen Ruan**
Indiana University

Huina Mao**
Indiana University

**Abstract**   Online social networking communities may exhibit highly complex and adaptive collective behaviors. Since emotions play such an important role in human decision making, how online networks modulate human collective mood states has become a matter of considerable interest. In spite of the increasing societal importance of online social networks, it is unknown whether assortative mixing of psychological states takes place in situations where social ties are mediated solely by online networking services in the absence of physical contact. Here, we show that the general happiness, or *subjective well-being* (SWB), of Twitter users, as measured from a 6-month record of their individual tweets, is indeed assortative across the Twitter social network. Our results imply that online social networks may be equally subject to the social mechanisms that cause assortative mixing in real social networks and that such assortative mixing takes place at the level of SWB. Given the increasing prevalence of online social networks, their propensity to connect users with similar levels of SWB may be an important factor in how positive and negative sentiments are maintained and spread through human society. Future research may focus on how event-specific mood states can propagate and influence user behavior in "real life."

## 1   Introduction

Bird flocking and fish schooling are typical and well-studied examples of how large communities of relatively simple individuals can exhibit highly complex and adaptive behaviors at the collective level. Agent-based models of such complex collective behaviors generally rest on the tenets that they are an emergent property of the local interactions between the individuals in a particular community and that these local interactions are governed by relatively simple sets of individual rules [14, 16, 19, 33, 39, 43].

An analogous situation may exist in the domain of social networking environments, where a plethora of complex and adaptive collective phenomena may emerge from the networked interactions between large groups of individuals who each perform relatively simple or minor tasks, such as answering questions, reporting the news, or updating their personal status. The recent role of social networking environments in social unrest [32] and large-scale disasters [38] illustrates that the collective behavior of online crowds may be of a complex psychosocial nature and exhibit some of the key features of complex adaptive systems [18], or in fact life [3].

---

 * Contact author.
** Indiana University, School of Informatics and Computing, 919 East 10th Street, Bloomington, IN 47408. E-mail: jbollen@indiana.edu (J.B.); bgoncalv@indiana.edu (B.G.); gruan@umail.iu.edu (G.R.); huinmao@umail.iu.edu (H.M.)

Emotions can play an important role in human decision making, and can thus shape collective phenomena that emerge in social networking environments. In fact, several studies have mentioned the emergence of collective online mood states that are highly adaptive and reactive to external events [5, 7, 24, 37, 41], and that in some cases have been shown to contain predictive information with regard to various socioeconomic indicators [6, 17]. A better understanding of how online social networking environments modulate human mood states is thus crucial to the development of models of how complex patterns of human collective behavior are shaped by their online interactions [26, 27, 47].

As the old adage goes, "Birds of a feather flock together." In network theory, this effect is known as *homophily* [29] or *assortative mixing* [34, 35, 36], and it occurs in a network when it disproportionately favors connections between vertices with similar characteristics. The opposite trend, that of favoring connections between nodes with different characteristics, is known as *disassortative mixing*. For example, a friendship network [29] may be highly assortative if it connects individuals who are at similar locations or have similar musical tastes. A heterosexual network [44], on the other hand, will be highly disassortative, since partners will tend to be of the opposite sex. However, few networks are entirely assortative or disassortative: Most will exhibit both properties to some degree depending on the particular characteristic.

Social networks can exhibit significant degrees of assortative mixing with respect to a variety of demographic attributes such as sex, race, age, religion, and education, including behavioral and health attributes [8, 11, 22, 31] and even genotypes [15]. Surprisingly, this also seems to be the case for certain psychological states such as loneliness [29]. Individuals preferentially share relations with individuals who report equally elevated levels of loneliness, and this homophilic tendency increases over time.

Although it is clear that psychological states affect behavior both online [10] and offline, the mechanisms through which such states exhibit assortativity and contagion across social bonds are not yet fully understood. However, two different processes are conceivable: that individuals seek homophilic social relations to share subjective experiences (homophilic attachment), or that the emotional state of an individual can influence that of the people with which he or she interacts (contagion) [42].

While both possibilities are clearly in play in real-world social interactions, it is not clear whether or not they are present in *online* social systems, which do not necessarily emerge from physical contact or in-person communication [21, 28]. The Twitter[1] microblogging service is a case in point. Twitter users can post brief personal updates of less than 140 characters at any time. These updates, known as *tweets*, are distributed to a limited group of *followers*, namely, other Twitter users who have elected to *follow* the particular user's tweets [23]. These follower relations are of a fundamentally different nature from their off-line counterparts [25]; they are not necessarily reciprocated (i.e., directed), nor modulated and are mostly focused on the exchange of information. In effect, a Twitter follower relation simply represents the fact that one individual is interested in the content produced by another, without the requirement that the interest be reciprocated. As a simple example, consider the case of celebrities that attract the attention and interest of a large number of people without reciprocating it. This arrangement results in a social network in the form of a directed, unweighted graph, which is quite different from naturally occurring social networks in which friendship ties are generally symmetric and vary in strength. As a consequence, one would expect homophily and assortative mixing of emotional states to be absent or fundamentally altered in online social networking environments, in particular those with asymmetric, unweighted connections such as Twitter.

However, in spite of the expectation that online environments fundamentally alter social interaction, recent results indicate that personal preferences do indeed exhibit homophilic properties in online environments such as BlogCatalog and Last.fm [4]. Tantalizingly, this has also been found the case for *sentiment* [49] in LiveJournal.[2]

Here we investigate whether and to what degree the general happiness, or *subjective well-being* (SWB) [12], of individual Twitter users exhibits homophilic properties or assortative mixing. Several

---

1 Twitter, http://www.twitter.com
2 LiveJournal: http://www.livejournal.com/

previous works have focused on aggregate [2, 5, 6, 13, 30] measurements of mood or emotion in entire communities or systems, but we analyze individual mood states in an online social network. On the basis of a collection of 129 million tweets, we track the SWB levels of 102,009 Twitter users over a 6-month period from the content of their tweets. Each is rated on an emotional scale using a standard sentiment analysis tool. A subsequent assortativity analysis of the Twitter social network then reveals its degree of SWB homophily or assortative mixing. Our results indicate that the overall SWB of Twitter users is positive, and highly assortative. In other words, Twitter users are preferentially connected to those with whom they share the same level of general happiness or SWB.

## 2  Data and Methods

We collected a large set of Tweets submitted to Twitter in the period from November 28, 2008 to May 2009. The data set consisted of 129 million tweets submitted by several million Twitter users. Each Tweet contained a unique identifier, date-time of submission (GMT+0), submission type, and textual content, among other information. Some examples are shown below in Table 1.

We complemented this cross-section sample of Twitter activity by retrieving the complete history of over 4 million users, as well as the identity of all of their followers. The final Twitter follower network contained 4,844,430 users (including followers of our users for which we did not collect timeline information). Armed with the social connections and activity of these users, we were able to measure the way in which the emotional content of each user varied in time and how it spread across links.

### 2.1  Creating a Twitter Friend Network

The follower network we collected consists of a directed graph $G = (V, E)$, in which $V$ is the set of all 4,844,430 Twitter users in our collection, and $E$ is the set of edges $E \subseteq V^2$, in which each directional edge $v \in E$ consist of a pair $(v_i, v_j)$ that indicates that user $v_i$ follows user $v_j$. By design, the Twitter social network is based on follower relations, which are unidirectional and very easy to establish. As such, they form a very minimal representation of possible interaction between those who follow and those who are being followed. In fact, it is quite common for a user $v_i$ to follow a user $v_j$, but for $v_j$ not to follow $v_i$ back. Thus, follower relations are not necessarily indicative of any personal relation that might in itself preclude the establishment of assortative mixing and homophily. We therefore distinguish between mere Twitter *followers* and actual *friends* [20, 46] by applying the following transformations to the Twitter follower graph $G$:

- First, we create a network of Twitter friend relations from the follower relations in $G$ by only retaining edges $(v_i, v_j) \in E$ for which we can find a *reciprocal* relation $(v_j, v_i)$. Thus the set of friend connections is $E' = \{(v_i, v_j) : \exists (v_j, v_i) \in E\}$; that is, two users only share a friendship tie if they are both following each other (Figure 1).

- Second, to exclude occasional users who are not truly involved in the Twitter social network, we only retained those users in our Twitter friend network who posted more than one tweet per day on average over the course of 6 months.

Table 1. Examples of tweet data collected from November 28, 2008 to May 2009 for 4,844,430 users.

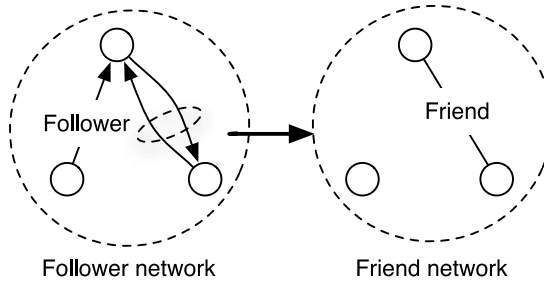| ID | Date-time | Type | Text |
|---|---|---|---|
| 1 | 2008-11-28 02:35:48 | Web | Getting ready for Black Friday. Sleeping out at Circuit City or Walmart not sure which. So cold out. |
| 2 | 2008-11-28 02:35:48 | Web | @dane I didn't know I had an uncle named Bob :-P I am going to be checking out the new Flip sometime soon |

Figure 1. Converting the original follower network of Twitter into a friend network by only taking into account reciprocal connections.

- Third, we assign a weight $w_{i,j}$ to each edge $(v_i, v_j)$, which serves as an indication of the degree to which users $v_i$ and $v_j$ have similar sets of friends:

$$w_{i,j} = \frac{\|C_i \cap C_j\|}{\|C_i \cup C_j\|} \tag{1}$$

  where $C_i$ denotes the neighborhood of friends surrounding user $v_i$. Note that this approach does not take into account the number of tweets exchanged between two users, but simply the degree to which two Twitter users have similar friends (Figure 2).

- Finally, we extracted the largest Connected Component ($G_{CC}$) from the resulting network, thereby obtaining a Twitter friend network of 102,009 users and 2,361,547 edges.

The reduction in nodes from our original Twitter follower network (4,844,430) to the resulting friend network (102,009) indicates that in Twitter only a small fraction of users are involved in the type of reciprocated follower type that we consider indicative of actual social relationships. However, once this reduction has occurred, we find that the largest connected component of the friend network, $G_{CC}$, retains 97.9% of users in the original Twitter friend network. This indicates a high degree of connectivity across all users in the final friend graph. The diameter of $G_{CC}$ was found to be only 14, indicating that in spite of the low density, two friends are separated by at most 14 steps. Other relevant network parameters for $G_{CC}$ are provided in Table 2 (Figure 3).

Examining the edge weight distribution as shown in Figure 4, we observe a strongly skewed frequency distribution, indicating very many connections in $G_{CC}$ with low edge weights ($w_{i,j} < 0.3$) and few connections with very high edge weights ($w_{i,j} > 0.6$). The degree frequency distribution reveals a similar pattern, with most users connected to only a few users and a small minority of users connected to thousands of users.
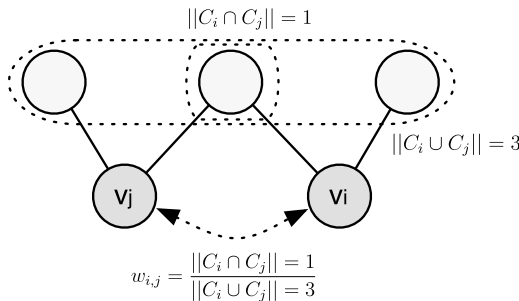


Figure 2. Example of Twitter friend similarity as calculated according to Equation 1. Users $v_i$ and $v_j$ share one friend out of three total. Therefore their connection is assigned a weight $w_{i,j} = \frac{1}{3}$.

Table 2. Network parameters for largest connected component of Twitter friend network.

| Network parameter | Value |
| --- | --- |
| Nodes | 102,009 users |
| Edges | 2,361,547 edges |
| Density | 0.000454 |
| Diameter | 14 |
| Average degree | 46.300 |
| Average clustering coefficient | 0.262 |

## 2.2   User-Level Measurements of Subjective Well-Being

We cannot directly interrogate Twitter users about their SWB [12], but we can infer users' SWB from the aggregate emotional content of their tweets over a period of 6 months. To do so we apply the following procedure.
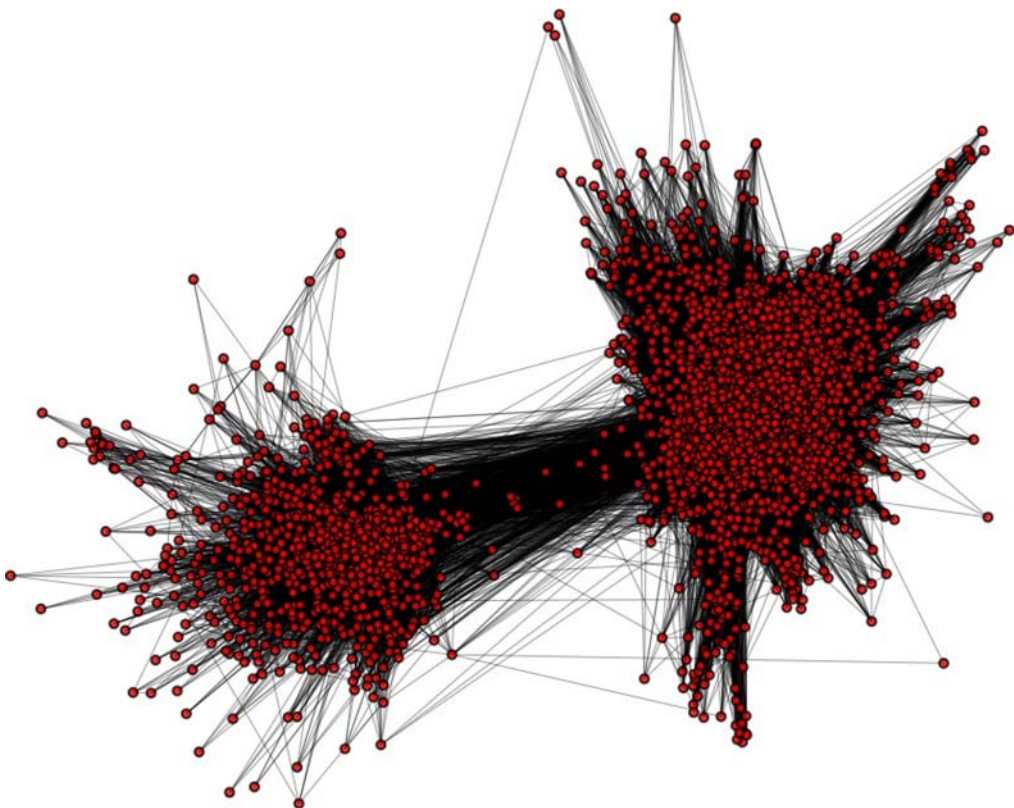


Figure 3. A subgraph of 3,587 users extracted from the generated Twitter social network (102,009 users and 2,361,547 edges), illustrating the high degree of clustering in the resulting network at the level of their friend relations (see Table 2). Clusters result from network connections. SWB is not visualized in this graph.
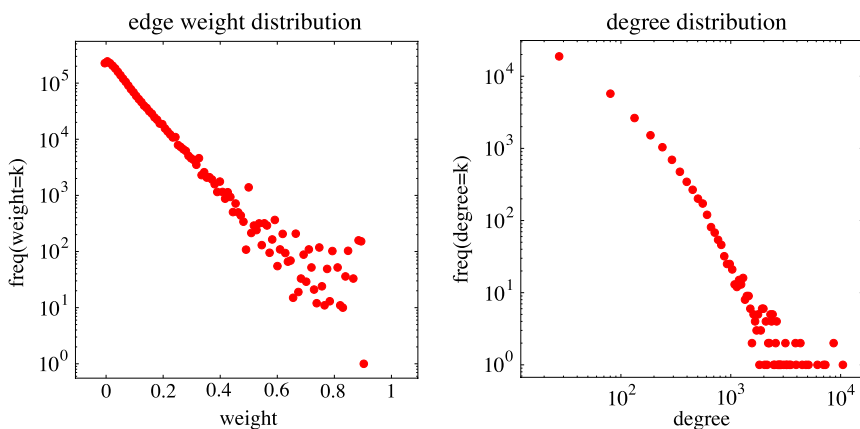
**Figure 4.** Twitter friend network edge weight and degree distributions.

To reduce noise, we only included Twitter users in $G_{CC}$ that posted at least one tweet per day. This guaranteed at least 180 tweets for every individual user from which to assess their SWB. We then analyzed the emotional content of each user's 6-month record of tweets, using OpinionFinder (OF),[3] which is a publicly available software package for sentiment analysis that can be applied to determine sentence-level subjectivity [48]. In other words, we use OF to identify the emotional polarity (positive or negative) of each tweet in a user's record. OF has been successfully used to analyze the emotional content of large collections of tweets [37] by using its lexicon to determine the dominance of positive or negative tweets on a given day. Here we select both positive and negative words that are marked as either "weak" and "strong" from the OF sentiment lexicon, resulting in a list of 2,718 positive and 4,912 negative words. For each tweet in an individual user's 6-month record, we count the numbers of negative and positive terms from the OF lexicon that it contains, and increase the individual user's score of either negative or positive tweets by 1 for each occurrence.

The subjective well-being $S(u)$ of user $u$ is then defined as the fractional difference between the number of tweets that contain positive OF terms and those that contain negative terms:

$$S(u) = \frac{N_p(u) - N_n(u)}{N_p(u) + N_n(u)}$$

where $N_p(u)$ and $N_n(u)$ represent, respectively, the numbers of positive and negative tweets for user $u$.

A number of examples are shown in Table 3.

## 2.3   Defining SWB Assortativity

Having calculated the SWB of each user, we can now proceed to measure the degree to which the SWBs of connected users are correlated. Intuitively, a person can be emotionally influenced by their friends in two complementary ways: Influence can come from interacting with specific individuals to which one may attribute more importance [40]. We refer to this first type as *pairwise node assortativity*, since it assesses the degree to which every two pairwise-connected users have similar SWB values. Another possibility is that each individual is influenced by the overall SWB of all of the people he or she interacts with. We refer to this second type as *neighborhood assortativity*.

---

3 http://www.cs.pitt.edu/mpqa/opinionfinderrelease/

Figure 5 illustrates this distinction; it shows the actual neighborhood friend network of an individual in $G_{CC}$ who has very high SWB values. Nodes are colored according to their SWB values, with red indicating high SWB values, blue indicating low SWB values, and white indicating neutral or zero SWB values. A particular individual with high SWB is connected to a local network of equally high-SWB individuals (red). The individual can thus be said to be neighborhood assortative within this cluster. However, the individual is also connected to several individuals with low SWB values (blue). For each individual connection this is a case of pairwise disassortativity. The cluster of low-SWB individuals, on the other hand, exhibits neighborhood assortativity for low SWB values, and the network in its entirety, including both low- and high-SWB clusters, exhibits strong SWB assortativity; nodes with similar low or high SWB values tend to be connected (blue and red clusters).

We formally define *pairwise* SWB assortativity as follows: For each edge $(v_i, v_j)$ in the $G_{CC}$ of our social network, we extract the corresponding two SWB values, one for the source node and one for the target node. These values are then aggregated into two vectors, $\mathcal{S}(S)$ and $\mathcal{S}(T)$ for sources and targets, respectively. The value of the pairwise assortativity, denoted $A_P(G_{CC})$, is then given by the Pearson correlation coefficient $\rho$ of these two vectors:

$$A_P(G_{CC}) \equiv \rho(\mathcal{S}(S), \mathcal{S}(T)) = \frac{1}{n-1} \sum_i \left[ \left( \frac{\mathcal{S}(S_i) - \langle \mathcal{S}(S) \rangle}{\sigma(\mathcal{S}(S))} \right) \left( \frac{\mathcal{S}(T_i) - \langle \mathcal{S}(T) \rangle}{\sigma(\mathcal{S}(T))} \right) \right] \tag{2}$$

The pairwise assortativity is thus defined in the interval $[-1, +1]$, with $-1$ indicating perfect disassortativity, 0 a lack of any assortativity, and $+1$ perfect assortativity.

The *neighborhood* assortativity of $G_{CC}$ with regard to SWB, denoted $A_N(G_{CC})$, can be calculated as follows. For each user $u \in V$, we define its neighborhood:

$$\kappa(u) = \{\forall v : \exists (u, v) \in E\} \tag{3}$$

**Table 3.** Examples of tweets posted by users with very high and very low SWB values.

| *Tweets submitted by high-SWB users (>0.5)* |
|---|
| So…nothing quite feels like a good shower, shave and haircut…love it |
| My beautiful friend. i love you sweet smile and your amazing soul |
| i am very happy. People in Chicago loved my conference. Love you, my sweet friends |
| @anonymous thanks for your follow I am following you back, great group amazing people |

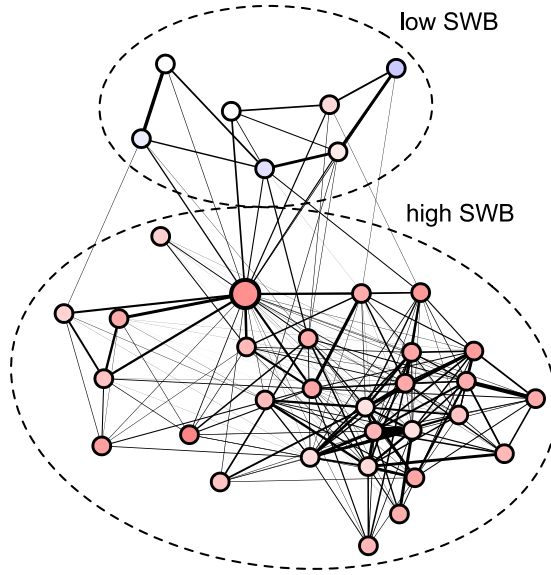| *Tweets submitted by low-SWB users (<0.0)* |
|---|
| She doesn't deserve the tears but i cry them anyway |
| I'm sick and my body decides to attack my face and make me break out!! WTF :( |
| I think my headphones are electrocuting me. |
| My mom almost killed me this morning. I don't know how much longer i can be here. |

**Figure 5.** Neighborhood network of a very high-SWB individual (center). Blue, white, and red node colors correspond respectively to low, neutral, and high SWB values.

so that $\kappa(u)$, or $\kappa_u$, represents the set of users that user $u$ is connected to. We then calculate an average SWB value for $\kappa(v)$, which we denote by

$$\overline{\mathcal{S}(\kappa_u)} = \frac{1}{\|\kappa(u)\|} \sum_{v \in \kappa(u)} \mathcal{S}(v) \tag{4}$$

We can now define two vectors, one for the SWB values of an individual user and one for the average SWB value of his or her neighborhood, denoted by $\mathcal{S}(U)$ and $\overline{\mathcal{S}(\kappa)}$. The neighborhood assortativity of the network $G_{CC}$ with regard to SWB, denoted $A_\kappa$, is then given by the correlation function $\rho$ computed over these two vectors as follows:

$$A_K(G_{CC}) \equiv \rho\left(\mathcal{S}(U), \overline{\mathcal{S}(\kappa)}\right) = \frac{1}{n-1} \sum_u \left[ \left(\frac{\mathcal{S}(u) - \langle \mathcal{S}(U) \rangle}{\sigma(\mathcal{S}(U))}\right) \left(\frac{\overline{\mathcal{S}(\kappa_u)} - \langle \overline{\mathcal{S}(\kappa)} \rangle}{\sigma\left(\overline{\mathcal{S}(\kappa)}\right)}\right) \right] \tag{5}$$

with the sum to be taken over every user $u$. Then $A_K(G_{CC})$ represents the correlation between the SWB values of user $v_i$ and the mean SWB values of his or her friends. This is similar to the pairwise version, $A_{\kappa_i}$ in the range $[-1, +1]$, where $-1$ indicates perfect neighborhood disassortativity, and where $+1$ indicates perfect neighborhood assortativity.

## 3  Results and Discussion

### 3.1  SWB Distribution

In Figure 6 we plot the probability distribution of SWB values across all Twitter users in our sample. The distribution seems bimodal with two peaks: one in the range $[-0.1, 0.1]$ and another in the range $[0.2, 0.4]$. Excluding users whose SWB $= 0$ (due to a lack of emotional content in their Tweets), we find that a majority of Twitter users in our sample have positive SWB values in a rather narrow range $[0.1, 0.4]$,
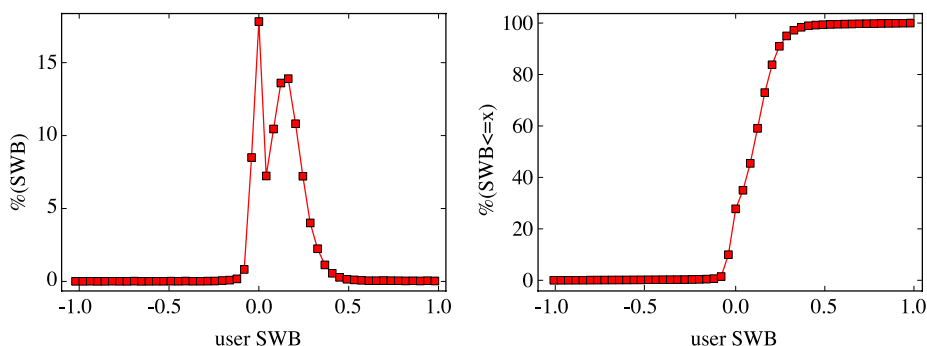
**Figure 6.** Probability density (%) and cumulative distribution (%) of subjective well-being (SWB) for our sample of Twitter users.

with a peak at SWB $= 0.16$. This is confirmed by the cumulative distribution shown at the bottom left of Figure 6; 50% of users have SWB values $\leq 0.1$, and 95% of users have SWB values $\leq 0.285$.

### 3.2 Pairwise and Neighborhood SWB Assortativity

In Equation 2 we defined the pairwise SWB as the correlation between the SWB values of connected users in our Twitter friend networks, whereas Equation 5 defined neighborhood assortativity as the correlation between the SWB values of individual users and the mean SWB values of their neighbors in the graph $G_{CC}$. The assortativity values were found to be[4] $0.443^{\star\star\star}$ ($N = 2{,}062{,}714$ edges) for the pairwise SWB assortativity, and $0.689^{\star\star\star}$ ($N = 102{,}009$) for the neighborhood assortativity. Both correlations are highly statistically significant ($p < .001$, indicated by $^{\star\star\star}$) for the sample sizes.

The scatterplots on the left of Figure 7 and Figure 8 show the distribution of SWB values[5] across the sample of all edges and nodes in $G_{CC}$ and confirm the observed correlation between the SWB values of connected or neighboring users in $G_{CC}$.

The pairwise assortativity scatterplot (Figure 7, left) indicates a significant amount of scatter, commensurate with the lower correlation value of $0.443^{\star\star\star}$, which is nevertheless statistically highly significant. The observed relation is not obviously linear. The distribution of values is affected by the bimodal distribution of SWB values as shown in Figure 6; large numbers of observations cluster at SWB values within the ranges $[-0.05, 0.05]$ and $[0.1, 0.3]$. The clustering pattern of Figure 7, however, indicates that users with SWB values in a particular range are preferentially connected to users within that same range, thereby confirming the observed positive pairwise SWB assortativity.

The neighborhood assortativity scatterplot (Figure 8, left) indicates a similar effect, but here the clustering of users is less pronounced and the amount of scatter is lower than that observed for the pairwise assortativity scatterplot, commensurate with the higher neighborhood assortativity value of $0.689^{\star\star\star}$. Although less pronounced, the bimodal distribution of SWB values is apparent and leads to a clustering of user and neighborhood SWB values in the ranges $[-0.05, 0.05]$ and $[0.1, 0.3]$. Nevertheless it is again the case that users with SWB values in either range are most likely to be connected to users or neighborhoods with SWB values in the same range. The distribution of user and neighborhood SWB value is furthermore in line with a positive linear relationship.

### 3.3 Edge Weight and SWB Assortativity

Pairwise SWB assortativity and neighborhood SWB assortativity diverge significantly ($0.443^{\star\star\star} < 0.689^{\star\star\star}$). The former is based on the pairwise comparison of SWB values across all connections

---

4 The sample sizes for pairwise assortativity and neighborhood assortativity are expressed in edges and nodes, respectively, since the former correlation is calculated on the basis of a sample of edges that connect pairs of nodes, whereas the latter is calculated on the basis of a sample of nodes and their neighborhoods.

5 Some banding occurs in the scatterplots, possibly due to common ratios of integers, for example 0.5, factoring into the SWB calculation.
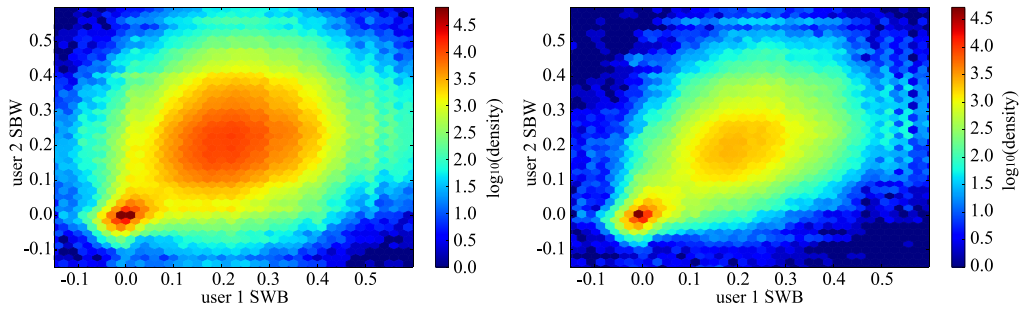
**Figure 7.** Scatterplot of SWB values for user connected in Twitter friend network. Left: All edges included. SWB assortativity = 0.443***, N = 2,062,714 edges. Right: Scatterplot includes edges $w_{i,j} \geq 0.1$, SWB assortativity = 0.712***, N = 479,401.

in $G_{CC}$, many of which may be weak or irrelevant connections from the perspective of indicating actual friend ties. To measure the effect of edge weights, we calculate pairwise and neighborhood assortativity values under different edge thresholds, that is, we only take into account edges in $G_{CC}$ whose weight as defined in Equation 1 is $w_{i,j} \geq \epsilon$, where $\epsilon \in [0, 1]$ represents a given edge threshold. The consequent assortativity calculations will therefore more strongly reflect only those connections between users that are indicative of stronger friend relations (higher $w_{i,j}$). In other words, we are verifying whether stronger user relations lead to higher or lower assortativity.

The results of the calculation of pairwise and neighborhood SWB assortativity under various edge thresholds are shown in Table 4 and visualized in Figure 9. Values for $\epsilon > 0.8$ are excluded, since the correlation coefficients were not statistically significant ($p < .1$). The graphs in Figure 9 overlay the different pairwise and neighborhood SWB assortativity values along with the number of remaining edges and nodes under the given edge threshold, that is, the sample size for the given assortativity calculation.

Both pairwise and neighborhood assortativity values increase as the edge threshold $\epsilon$ increases, but not in a linear manner. Pairwise SWB assortativity values increase sharply as $\epsilon$ increases from 0 to 0.10 and afterward stabilize at a value of approximately 0.750, which is maintained in the interval $\epsilon \in [0.15, 0.85]$. In other words, removing edges with $w_{i,j} < 0.1$ increases pairwise SWB assortativity considerably, but the removal of edges with higher $w_{i,j}$ values has little or no additional effect. We observe that $\epsilon = 0.1$ reduces the number of edges by a fifth—namely, from 2,062,714 to 479,401— indicating that a large number of edges are characterized by low similarity values, which when included lead to lower pairwise SWB assortativity.

The neighborhood SWB assortativity increases with higher $\epsilon$ values, but in a less pronounced manner. At $\epsilon = 0$ we find a neighborhood SWB assortativity value of 0.689, which increases to approximately 0.760 for $\epsilon \in [0.10, 0.90]$. We observe significant declines in the number of nodes that remain under increasing $\epsilon$ values, as was the case for pairwise assortativity.
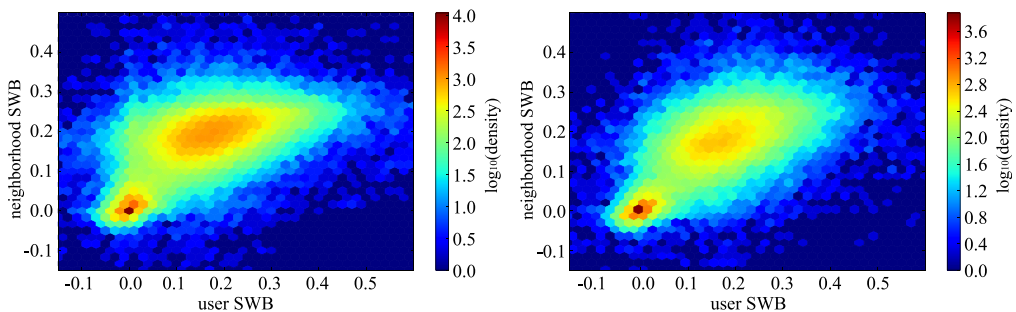


**Figure 8.** Scatterplot of SWB values for users (*x*) and their neighborhood (*y*) connected in Twitter friend network. Left: All edges included. SWB assortativity = 0.689***, N = 102,009 nodes. Right: Scatterplot includes edges $w_{i,j} \geq 0.1$, SWB assortativity = 0.746***, N = 59,952.

**Table 4.** Pairwise and neighborhood subjective well-being assortativity values *A*(SWB) versus edge threshold ∈. ***: *p* < .001.

| Edge threshold ∈ | Pairwise | | Neighborhood | |
|---|---|---|---|---|
| | *A*(SWB) | *N* edges | *A*(SWB) | *N* nodes |
| 0.0 | 0.443*** | 2,062,714 | 0.689*** | 102,009 |
| 0.10 | 0.712*** | 479,401 | 0.746*** | 59,952 |
| 0.20 | 0.754*** | 128,261 | 0.769*** | 33,693 |
| 0.30 | 0.755*** | 36,255 | 0.780*** | 16,334 |
| 0.40 | 0.743*** | 10,355 | 0.779*** | 7,699 |
| 0.50 | 0.757*** | 3,255 | 0.781*** | 3,793 |
| 0.60 | 0.798*** | 1,375 | 0.805*** | 1,439 |
| 0.70 | 0.755*** | 689 | 0.816*** | 502 |
| 0.80 | 0.434*** | 301 | 0.768*** | 149 |
| 0.90 | — | — | — | — |

At a value of ∈ = 0.1 we find the highest pairwise and neighborhood SWB assortativity values combined with the largest sample sizes, excluding the absence of threshold (i.e., ∈ = 0). We therefore generated new scatterplots of SWB values for pairwise and neighborhood SWB assortativity at ∈ = 0.1, as shown in Figure 7 (right) and Figure 8 (right). The scatterplots reflect higher assortativity values; we find less scatter, a stronger positive and linear relation between SWB value of connected users, and a less pronounced clustering caused by the bimodal distribution of SWB values.

### 3.4 Discussion

The above outlined results indicate the following.

First, the Twitter users in our sample exhibit moderate to high SWB, with very few being characterized by very low SWB values. One may thus conclude that they are a moderately happy group.
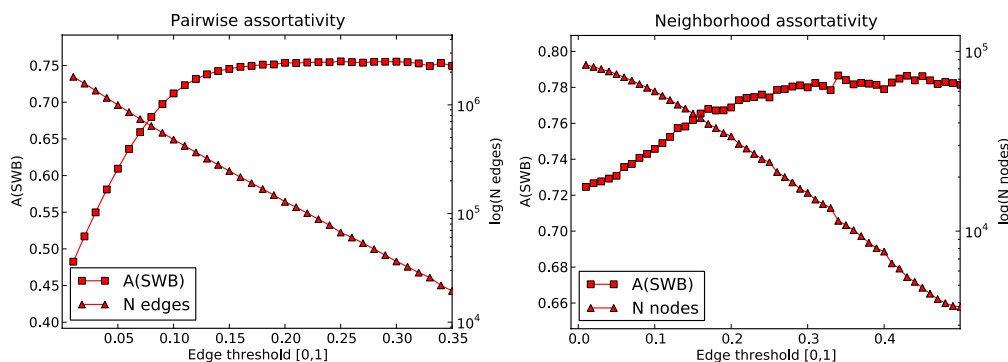


**Figure 9.** Subjective well-being assortativity and log(number of edges) versus edge weight thresholds for pairwise and neighborhood assortativity.

This observation is most likely an underestimation given the relative preponderance of negative terms in the OF lexicon. However, the SWB distribution is bimodal, showing a group of Twitter users with zero to very low SWB values (i.e., those that are on the average somewhat happy) and another group with more pronounced, higher SWB values. This may result from sociocultural differences in how emotions and mood are expressed on Twitter. Some users may express their emotional states infrequently, whereas others are more prone to do so.

Second, we find statistically significant levels of pairwise and neighborhood SWB assortativity, indicating that Twitter users either prefer the company of users with similar SWB values (assumption) or converge on their friends' SWB values (contagion). The relation between user SWB values is not linear and is biased by the bimodal distribution of SWB values, users being clustered in two groups with equally low or equally high SWB values. In other words, low-SWB users are connected to low-SWB users, and high-SWB users are connected to high-SWB users. Again, this may confirm the notion that distinct sociocultural factors affect the expression of emotion and mood on Twitter, and cause users to cluster according to their degree of expressiveness as well as SWB.

The results of measuring pairwise and neighborhood assortativity under different edge weights indicate that we find a stronger and more significant relation between the SWB values of connected users when we only take into account connections with higher weights $w_{i,j}$, that is, those that are deemed more reliable indicators of actual friend ties. A possible mechanism to explain this effect might be that users' neighborhoods contain individuals that they are indeed strongly assortative with, and whose SWB values affect mean neighborhood SWB values and thus neighborhood assortativity overall, whereas they are drowned out in the process of making pairwise comparisons between all individuals that a user is connected to in the process of calculating pairwise SWB assortativity. For example, users might generally have 10 neighbors, but be generally highly SWB assortative with only one friend. In the calculation of pairwise assortativity, this leads to 10 pairwise comparisons between SWB values, only one of which contributes to the overall observed pairwise SWB assortativity in the graph. However, the neighborhood assortativity relies on a mean SWB value calculated for the entire neighborhood, including the single highly assortative individual. The latter thus influences the average SWB value for the entire neighborhood, causing an increased neighborhood assortativity value.

The greatest improvement in assortativity values indeed occurs for pairwise SWB assortativity, which is most affected by the preponderance of weakly weighted connections, since it is defined at the level of all individual user-to-user connections. Both pairwise and neighborhood SWB assortativity converge on a value of approximately 0.750, which indicates a significant degree of SWB assortativity in our Twitter friend graph $G_{CC}$.

## 4   Conclusion

Recent findings show that assortative mixing can occur in a variety of social contexts and personal attributes. Here we show that subjective well-being is equally assortative in the Twitter social network, that is, the SWB of individuals that have reciprocal Twitter follower links are strongly related. Happy users tend to connect to happy users, whereas unhappy users tend to be predominantly connected to unhappy users. The convergence of pairwise and neighborhood assortativity under increasing edge weight thresholds indicates that users tend to be most assortative with a limited number of individuals that they have strong social ties to and that weaker ties fulfill a different social role, possibly as outlined by Centola and Macy [9].

We do not not address the social or cognitive mechanisms that cause the observed SWB assortativity. Two different mechanisms may be at work [1]. The first is based on the notion of *homophily*, that is, users and connections tend to preferentially connect to users with similar SWB values. As an online social network grows, new connections are thus biased toward connecting individuals with similar SWB values. This process may be modeled in terms of *preferential attachment* theory. The second mechanism that may cause SWB assortativity is that of *mood contagion*, namely, that connected users converge to similar SWB values over time. In other words, being connected to unhappy users

can make one unhappier, and vice versa. The latter suggests that users may control their own level of SWB by choosing the right set of online friends and influence their friends' SWB by creating strong social ties and hoping for some form of SWB contagion to take place. A third possibility is that users assess or express their SWB relative to that of their friends. As a user's neighborhood becomes happier, this may affect their own expression of SWB-related sentiment. This phenomenon may occur at the level of entire cultures, which may be more prone or less prone to open expressions of individual sentiment.

At this point our research does not offer any information on which of these mechanisms causes the observed SWB assortativity, or in fact whether both may be occurring. Future research will therefore focus on analyzing user connections and SWB values over time, and relating these changes in the framework of homophily and preferential attachment [45]. Twitter has now become a major international phenomenon, and this investigation must therefore include linguistic, cultural, and geographic factors.

## Acknowledgments

## References

1. Aral, S., Muchnik, L., & Sundararajan, A. (2009). Distinguishing influence based contagion from homophily-driven diffusion in dynamic networks. *Proceedings of the National Academy of Sciences of the United States of America*, *106*, 21544–21549.

2. Balog, K., Mishne, G., & de Rijke, M. (2006). Why are they excited? Identifying and explaining spikes in blog mood levels. In *Proceedings of the 11th Conference of the European Chapter of the Association for Computational Linguistics*.

3. Bedau, M. A., McCaskill, J. S., Packard, N. H., & Rasmussen, S. (2010). Living technology: Exploiting life's principles in technology. *Artificial Life*, *16*, 89–97.

4. Bisgin, H., Agarwal, N., & Xu, X. (2010). Investigating homophily in online social networks. In *2010 IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology* (pp. 533–536).

5. Bollen, J., Mao, H., & Pepe, A. (2010). Determining the public mood state by analysis of microblogging posts. In *Proceedings of the Alife XII Conference* (pp. 667–669).

6. Bollen, J., Mao, H., & Zeng, X. (2011). Twitter mood predicts the stock market. *Journal of Computational Science*, *2*(1), 1–8. doi:10.1016/j.jocs.2010.12.007, arXiv: abs/1010.3003.

7. Bollen, J., Pepe, A., & Mao, H. (2011). Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena. http://arxiv.org/abs/0911.1583. Presented at the Fifth International AAAI Conference on Weblogs and Social Media (ICWSM11).

8. Cacioppo, J. T., Fowler, J. H., & Christakis, N. A. (2010). Alone in the crowd: The structure and spread of loneliness in a large social network. *Journal of Personality and Social Psychology*, *97*, 977–991.

9. Centola, D., & Macy, M. (2007). Complex contagion and the weakness of long ties. *American Journal of Sociology*, *113*, 702–734.

10. Chmiel, A., Sienkiewicz, J., Paltoglou, G., Buckley, K., Thelwall, M., & Holyst, J. A. (2010). Negative emotions boost users activity at BBC forum (Technical Report 1011.5459). arXiv.

11. Christakis, N. A., & Fowler, J. H. (2007). The spread of obesity in a large social network over 32 years. *The New England Journal of Medicine*, *357*, 370–379.

12. Diener, E. (2009). *Subjective well-being* (pp. 11–58). Berlin: Springer.

13. Dodds, P. S., & Danforth, C. M. (2009). Measuring the happiness of large-scale written expression: Songs, blogs, and presidents. *Journal of Happiness*, *11*.

14. Dorigo, M., Bonabeau, E., & Theraulaz, G. (2000). Ant algorithms and stigmergy. *Future Generation Computer Systems*, *16*, 851–871.

15. Fowler, J. H., Settle, J. E., & Christakis, N. A. (2011). Correlated genotypes in friendship networks. *Proceedings of the National Academy of Sciences of the USA*, *108*, 1993.

16. Garnier, S., Jost, C., Gautrais, J., Asadpour, M., Caprari, G., Jeanson, R., Grimal, A., & Theraulaz, G. (2008). The embodiment of cockroach aggregation behavior in a group of micro-robots. *Artificial Life*, *14*, 387–408.

17. Gilbert, E., & Karahalios, K. (2010). Widespread worry and the stock market. In *Fourth International AAAI Conference on Weblogs and Social Media* (pp. 58–65).

18. Holland, J. H. (1992). Complex adaptive systems. *Daedalus*, *121*, 17–30.

19. Holland, O. E., & Melhuish, C. (1999). Stigmergy, self-organization, and sorting in collective robotics. *Artificial Life*, *5*(2), 173–202.

20. Huberman, B. A., Romero, D. M., & Wu, F. (2008). Social networks that matter: Twitter under the microscope. *First Monday*, *14*, 1.

21. Huntsinger, J. R., Lun, J., Sinclair, S., & Clore, G. L. (2009). Contagion without contact: Anticipatory mood matching in response to affiliative motivation. *Personality and Social Psychology Bulletin*, *35*, 909–922.

22. Ibarra, H. (1992). Homophily and differential returns: Sex differences in network structure and access in an advertising firm. *Administrative Science Quarterly*, *37*, 422.

23. Java, A., Song, X., Finin, T., & Tseng, B. (2007). Why we twitter: Understanding microblogging usage and communities. In *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis* (pp. 56–65).

24. Kramer, A. D. I. (2010). An unobtrusive behavioral model of "gross national happiness." In *ACM CHI 2010: Conference on Human Factors in Computing* (pp. 287–290).

25. Kwak, H., Lee, C., Park, H., & Moon, S. (2010). What is Twitter, a social network or a news media? In *WWW '10: Proceedings of the 19th International Conference on World Wide Web* (pp. 591–600).

26. Lazer, D., et al. (2009). Computational social science. *Science*, *323*(5915), 721–723.

27. Lee, S. H., Kim, P.-J., Ahn, Y.-Y., & Jeong, H. (2010). Googling social interactions: Web search engine based social network construction. *PLoS One*, *5*, e11233.

28. Mason, W. A., Conrey, F. R., & Smith, E. R. (2007). Situating social influence processes: Dynamic, multidirectional flows of influence within social networks. *Personality and Social Psychology Review*, *11*, 279–300.

29. McPherson, M., Smith-Lovin, L., & Cook, J. M. (2001). Birds of a feather: Homophily in social networks. *Annual Review of Sociology*, *27*, 415–444.

30. Mislove, A., Lehmann, S., Ahn, Y.-Y., Onnela, J.-P., & Rosenquist, J. N. (2010). Pulse of the nation: Visualizing the mood of twitter. http://www.ccs.neu.edu/home/amislove/twittermood/

31. Mollica, K. A., Gray, B., & Treviño, L. K. (2003). Racial homophily and its persistence in newcomers' social networks. *Organization Science*, *14*, 123–136.

32. Mungiu-Pippidi, A., & Muntean, I. (2009). Moldova's "Twitter" revolution. *Journal of Democracy*, *20*.

33. Nathan, A., & Barbosa, V. C. (2008). V-like formations in flocks of artificial birds. *Artificial Life*, *14*, 179–188.

34. Newman, M. E. J. (2002). Assortative mixing in networks. *Physical Review Letters*, *89*, 208701/1–4.

35. Newman, M. E. J. (2003). Mixing patterns in networks. *Physical Review E*, *67*, 26126.

36. Newman, M. E. J. (2003). The structure and function of complex networks. *SIAM Review*, *45*, 167.

37. O'Connor, B., Balasubramanyan, R., Routledge, B. R., & Smith, N. A. (2010). From tweets to polls: Linking text sentiment to public opinion time series. In *Proceedings of the International AAAI Conference on Weblogs and Social Media*.

38. Oh, O., Kwon, K. H., & Rao, H. R. (2010). An exploration of social media in extreme events: Rumor theory and Twitter during the Haiti earthquake. In *Proceedings of the International Conference on Information Systems (ICIS)* (paper 231).

39. Olfati-Saber, R. (2006). Flocking for multi-agent dynamic systems: Algorithms and theory. *Engineering*, *51*, 401–420.

40. Onnela, J.-P., & Reed-Tsochas, F. (2010). Spontaneous emergence of social influence in online systems. *Proceedings of the National Academy of Sciences of the USA, 107*, 18375.

41. Pak, A., & Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. In N. Calzolari, K. Choukri, B. Maegaard, J. Mariani, J. Odijk, S. Piperidis, M. Rosner, & D. Tapias (Eds.), *Proceedings of the Seventh Conference on International Language Resources and Evaluation (LREC'10)*.

42. Parkinson, B., & Simons, G. (2009). Affecting others: Social appraisal and emotion contagion in everyday decision making. *Personality and Social Psychology Bulletin, 35*, 1071–1084.

43. Reynolds, C. W. (1987). Flocks, herds, and schools: A distributed behavioral model. *Computer Graphics, 21*, 25–34.

44. Rocha, L. E. C., Liljeros, F., & Holme, P. (2010). Information dynamics shape the sexual networks of internet-mediated prostitution. *Proceedings of the National Academy of Sciences of the USA, 107*, 5706.

45. Shalizi, C. R., & Thomas, A. C. (2011). Homophily and contagion are generically confounded in observational social network studies. In *Sociological Methods and Research*. arXiv:1004.4704v3.

46. Szell, M., Lambiotte, R., & Thurner, S. (2010). Multirelational organization of large-scale social networks in an online world. *Proceedings of the National Academy of Sciences of the USA, 107*.

47. Vespignani, A. (2009). Predicting the behavior of techno-social systems. *Science, 325*, 425–428.

48. Wilson, T., Wiebe, J., & Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing—HLT05* (pp. 347–354).

49. Zafarani, R., Cole, W. D., & Liu, H. (2010). Sentiment propagation in social networks: A case study in LiveJournal. In S.-K. Chai, J. Salerno, & P. Mabry (Eds.), *SBP 2010—Advances in social computing* (pp. 413–420). Berlin: Springer.