

# Temporal Causality of Social Support in an Online Community for Cancer Survivors

Ngot Bui, John Yen, and Vasant Honavar

College of Information Sciences and Technology  
The Pennsylvania State University  
University Park, PA 16802, USA  
{npb123, jyen, vhonavar}@ist.psu.edu

**Abstract.** Online health communities (OHCs) constitute a useful source of information and social support for patients. American Cancer Society’s Cancer Survivor Network (CSN), a 173,000-member community, is the largest online network for cancer patients, survivors, and caregivers. A discussion thread in CSN is often initiated by a cancer survivor seeking support from other members of CSN. It captures a multi-party conversation that often serves the function of providing social support e.g., by bringing about a change of sentiment from negative to positive on the part of the thread originator. While previous studies regarding cancer survivors have shown that members of OHC derive benefits from their participation in such communities, causal accounts of the factors that contribute to the observed benefits have been lacking. This paper reports results of a study that seeks to address this gap by discovering temporal causality of the dynamics of sentiment change (on the part of the thread originators) in CSN. The resulting accounts offer new insights that the designers, managers and moderators of an online community such as CSN can utilize to facilitate and enhance the interactions so as to better meet the social support needs of the community participants. The proposed methodology also has broad applications in the discovery of temporal causality from big data.

## 1 Introduction

World Health Organization [1], estimated that 14.1 million new cancer cases and 8.2 million cancer-related deaths occurred worldwide in 2012. In 2014, the number of deaths due to cancer in the US was estimated to be in excess of 0.58 million and the number of new cancer cases diagnosed was estimated to be 1.66 million [2]. According to National Cancer Institute, approximately 13.7 million Americans with a history of cancer were alive on January 1, 2012 [2]. Some of these individuals were cancer free, while others still showed cancer symptoms and may have been under treatment [2].

About 72% of Internet users in the U.S. utilize the Internet for health-related purposes and 26% have read or watched someone else’s experience about health or medical issues during the previous year [3]. Hence, many cancer survivors rely on an online health community (OHC) for social support. Social support can help cancer survivors cope better with their condition and hence improve the quality of their lives [4]. A cancer OHC that includes both survivors and their caregivers provides a forum to share experiences about their cancer, cancer treatment, and daily living issues. Through such

online interaction, they support one another in ways that family members, friends or even health care providers often cannot [5]. Several studies have documented the benefits derived by cancer survivors through participation in an OHC. The benefits of OHC participation include increased social support [6, 4], reduced levels of depression, stress, and psychological trauma [7], increased optimism about their life with cancer [6], increased ability to cope with their disease, and improvements on the physical and the mental aspects of their lives [4, 8]. Several studies have attempted to quantify the benefits of OHCs using sentiment analysis [9, 10]. Sentiment analysis and topic modeling were applied to CSN breast and colorectal cancer forums to investigate how change in sentiment of thread originators (persons who start the thread) varies by discussion topic [9]. Furthermore, sentiment analysis has been used to classify posts in an online health forum to determine whether a thread in the forum could benefit from the moderator’s intervention [11]. Sentiment classification of user posts in an online cancer support community [12] has been used to determine the utility of social support, and information support to patients [13].

However, none of the preceding studies have provided a causal account of sentiment change on the part of OHC participants. Such an account could help the designers, managers and moderators of an online community such as CSN to facilitate and enhance the interactions so as to better meet the social support needs of the community participants. Against this background, we introduce a novel approach to uncover the temporal causality of the dynamics of sentiment change (on the part of the thread originators) in OHCs. The proposed approach leverages the machinery of temporal causality introduced in [14, 15] to analyze the temporal ordering of sentiments of participants (i.e., thread originator and repliers) in a thread from CSN forum for the causes of the final sentiments (either positive or negative) of the thread originator.<sup>1</sup>

The rest of the paper is organized as follows. Section 2 introduces the key notions of temporal causality and the machinery for reasoning about causes of events. Section 3 describes our approach in temporal causality analyzing in CSN. Section 4 presents results of experiments that demonstrate the utility of the proposed approach. Section 5 concludes with a summary and a discussion and an outline of some promising directions for further research.

## 2 Temporal Causality

We start by reviewing a few key notions [17]: An event  $B$  is said to be a *prima facie* or potential cause of an event  $A$  iff (i)  $B$  precedes  $A$ ; (ii) the probability of  $B$ ,  $P(B) > 0$ ; and (iii) the conditional probability  $P(A|B) > P(A)$ . We say that  $B$  is a spurious cause of  $A$  iff  $B$  is a *prima facie* cause of  $A$ , and there is an event  $C$  that precedes  $B$  such that (i)  $P(B, C) > 0$ ; (ii)  $P(A|B, C) = P(A|C)$ ; and (iii)  $P(A|B, C) \geq P(A|B)$ . That is,  $C$  occurs before  $B$  and accounts for the effect  $A$  as well as  $B$  does. For example, assume that smoking and yellow finger precede the development of lung cancer and both appear

<sup>1</sup> Unlike the approach introduced in [16], this approach explicitly captures the temporality of the relationship between cause and effect. In addition to being able to represent properties being true for durations of time, this also allows a direct representation of the time window between cause and effect.

to be the causes to lung cancer. However, yellow finger and lung cancer have the same common cause (i.e., smoking). A *prima facie* cause that is not a spurious cause is said to be a genuine cause [17]. Suppes [17] offers a method for testing whether a cause is spurious in the restricted setting where there are only two possible causes. Kleinberg [14] argued for a more stringent criterion, for a *prima facie* cause to be considered a genuine cause and introduced a method for assessing the causal significance of a potential cause of an effect which can be used to identify a genuine cause of an event from among a set of its potential causes.

Kleinberg's framework uses temporal logic [18] to represent and reason about events that occur in time. Temporal logic is a variant of propositional modal logic that admits the truth value of a formula, constructed from atomic propositions (sentences which are either true or false and encoded by propositional symbols) using logical connectives, (i.e., conjunction, disjunction, and negation) to be time dependent. Hence, temporal logic can be used to represent whether a property is true at some specific time. Computation Tree Logic (CTL) [19], a branching-time logic, can be used to represent the fact that a property will be true at some time in the future (e.g., at some point in the future, the train will arrive). Probabilistic Computation Tree Logic (PCTL) [20] extends CTL by specifying deadlines (requiring a property to hold before a specified window of time elapses) and quantifying the transition probability between the states in CTL. Probabilistic Kripke structures [19] can be used to represent and reason in PTCL.

**Definition 1. Probabilistic Kripke structure** Let  $AP$  be a set of atomic propositions, a probabilistic Kripke structure is a four tuple:  $K = \langle S, s^i, L, \mathcal{T} \rangle$ , where  $S$  is a finite set of states;  $s^i \in S$  is an initial state;  $L : S \rightarrow 2^{AP}$  is a labeling function assigning subsets of  $AP$  to states; and  $\mathcal{T} : S \times S \rightarrow [0, 1]$  is a transition probability function such that  $\forall s \in S : \sum_{s' \in S} \mathcal{T}(s, s') = 1$ .

To represent the time and probability in PTCL formulas, Hansson et. al. [20] provide three types of *modal operators* which are path operators:  $A$  ("for all path") and  $E$  ("for some future path"), temporal operators:  $G$  ("holds for entire future path") and  $F$  ("finally holds"), and the "lead-to" operator. The lead-to operator, which is useful in formalizing temporal priority for causality, is defined as  $f \rightsquigarrow_{\geq p}^{\leq t} g \equiv AG \left[ f \rightarrow F_{\geq p}^{\leq t} g \right]$  which means that whenever  $f$  holds there is a probability of at least  $p$  that  $g$  will hold via some series of transitions taking less than or equal  $t$  time units. Some scenarios require the specification of a lower bound of time for  $g$  to hold. In this case, the lead-to operator can be constructed as  $f \rightsquigarrow_{\geq p}^{\geq r, \leq s} g$ , which denotes that  $g$  must hold in between  $r$  and  $s$  time units with probability  $p$  where  $0 \leq r \leq s \leq \infty$  and  $s \neq \infty$ .

Armed with the machinery of PTCL, we can define a *prima facie* (or a potential) cause as follows [21]:

**Definition 2. Prima Facie Cause in PCTL**  $c$  is a *prima facie* cause of  $e$  if there is a  $p$  such that all three following conditions hold. (i)  $F_{>0}^{\leq \infty} c$ ; (ii)  $c \rightsquigarrow_{\geq p}^{\geq 1, \leq \infty} e$ , and (iii)  $F_{< p}^{\leq \infty} e$  where  $c$  and  $e$  are PTCL formulas.

For  $c$  to be a *prima facie* cause of an effect  $e$ , the state where  $c$  is true should be reachable with non-zero probability, and the probability of reaching a state where  $e$  is true from a

state where  $c$  is true should be greater than the probability of reaching a state where  $e$  is true from the initial state of the system. This can be interpreted as requiring that  $c$  must occur at least once, and that the conditional probability of  $e$  given  $c$  is greater than the marginal probability of  $e$ .

We adopt the technique from [20] to calculate the probabilities  $F^{\leq\infty}c$ ,  $c \rightsquigarrow^{\geq 1, \leq\infty} e$ , and  $F^{\leq\infty}e$  where  $F^{\leq\infty}e$  denotes the path probabilities summed over the set of all paths starting from the initial state of the  $K$  structure and ending at a state where  $e$  is true (means that probability of  $e$  regardless of  $c$ );  $c \rightsquigarrow^{\geq 1, \leq\infty} e$  denotes the path probabilities summed over the set of all paths starting from the state where  $c$  is true and ending at a state where  $e$  is true (i.e., the probability of  $e$  given  $c$ ).

### 3 Methodology

#### 3.1 Learning Sentiment Classifier for Posts

Since we cannot detect the sentiment of a CSN member directly, we use the sentiment expressed in a post as a proxy for the sentiment of the CSN member at the time the post was created. As in [22], we categorize the posts as expressing either a positive sentiment or a negative sentiment.

Table 1: Features of a Post

Feature Name	Description
PosLength	The number of words in the post
PosStrength	Positive Sentiment Strength of the post
NegStrength	Negative Sentiment Strength of the post
Neg	NumberOfNeg / PostLength, where NumberOfNeg is equal to number of negative words and emotions
PosVsNeg	(NumberOfPos + 1) / (NumberOfNeg + 1), where NumberOfPos is equal to number of positive words and emotions
Name	NumberOfName / PostLength, where NumberOfName is the number of names mentioned in the post
Slang	Number of Internet slang words in the post

We used a set of 298 posts which are manually classified by two independent annotators into one of two categories: positive or negative. As in [10], we extract seven features (see Table 1) from a post to train a predictor for assigning posts to the positive or negative category. SentiStrength [22] is used to extract *PosStrength* and *NegStrength* which represent for the positive sentiment strength and negative sentiment strength of the post, respectively. We make use of the four lists<sup>2</sup>, a list of positive and negative words [23], a list positive and negative emotions<sup>3</sup>, a list of popular English female and male names, and a list of Internet slang words to calculate the *Neg*, *PosVsNeg*, *Name*, and *Slang* features. We use Adaboost, which has been shown to generate the best sentiment classification model [10], to classify the posts. The Adaboost sentiment classifier

<sup>2</sup> <http://sites.google.com/site/qiubaojun/psu-sentiment.zip>

<sup>3</sup> [http://en.wikipedia.org/wiki/List\\_of\\_emotions](http://en.wikipedia.org/wiki/List_of_emotions)

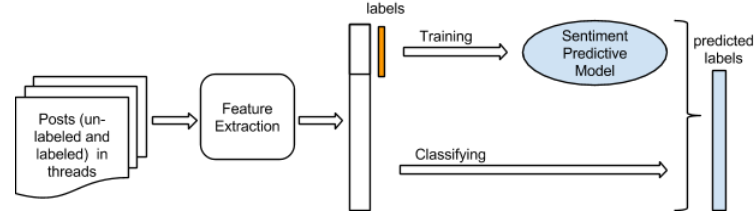


Fig. 1: Learning Sentiment Predictive Model

outputs a probability that a post expresses a positive sentiment, i.e.,  $\Pr(\text{positive} | \text{post})$ . If  $\Pr(\text{positive} | \text{post}) > 0.5$ , the post is classified as positive; otherwise, it is classified as negative. Figure 1 shows the procedure for training the post sentiment predictor and using it to classify new posts.

### 3.2 Cancer Survivor Network Thread as a Sequence of Sentiments

Cancer Survivor Network of American Cancer Society is an OHC with over 173,000 registered members which include cancer patients, their friends and families, and informal caregivers. In this paper, we use the CSN data set that contains all threads initiated between July 2000 and October 2010. The data set contains 48,779 discussion threads and more than 468,000 posts from 27,173 users. The data set is appropriately anonymized to protect the privacy of the CSN members.

Our goal is to uncover the causal effect (if any) of the temporally ordered posts that make up the thread on the final sentiment of the thread originators. More specifically, we are interested in discovering causal relationship between the reply posts and the change of sentiment of those who initiate the thread. Therefore, threads used in this study need to have at least one reply and at least one self-reply (i.e., a post by the thread originator later on thread). As a result, threads that don't contain a self-reply or reply are removed from this study. The resulting data set consists of 22,854 threads.

A thread can be represented as a temporally ordered sequence of posts  $P_{o1}, P_{r1}, P_{r2}, \dots, P_{o2}, \dots, P_{rm}, P_{on}$  where  $P_{o1}$  is the initial post from the thread originator;  $P_{oi}$  ( $i > 1$ ) are self-replies; and  $P_{rj}$  are replies to the post by individuals other than the thread originator. Since we focus on the communication between two kinds of actors in a thread, the thread originator and the individuals (other than the originator) who respond to the originator's post, we simply compute the average probability of positive sentiment of replies between two consecutive self-replies and use it as the positive sentiment probability of the collection of replies. Formally, the average positive sentiment probability is calculated as

$$\bar{p}_{ri} = \frac{\sum \Pr(\text{positive} | P_{rj})}{|P_{rj}|}, \forall j : t_{oi} \leq t_{rj} \leq t_{o(i+1)} \quad (1)$$

where  $t_{oi}$ ,  $t_{o(i+1)}$ , and  $t_{rj}$  are time points when posts  $P_{oi}$  and  $P_{o(i+1)}$  and  $P_{rj}$  are created, respectively and  $|P_{rj}|$  is the number of reply posts from  $t_{oi}$  to  $t_{o(i+1)}$ .

To establish a formal representation that enables temporal causality analysis, we transform the sequence of post sentiment probability in a thread to a sequence of post

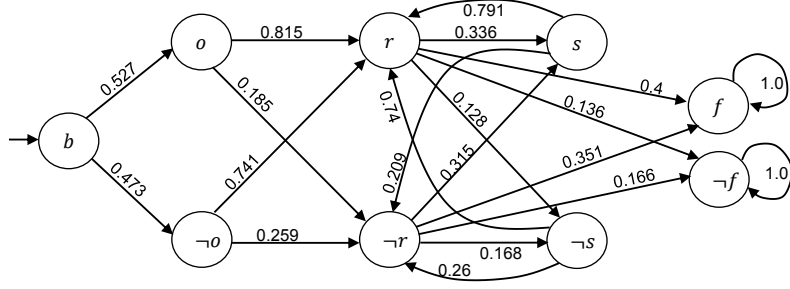


Fig. 2: Probabilistic Kripke Structure for CSN

sentiment states as follows:  $[Sentiment\ state\ of\ initial\ post]\ [Average\ sentiment\ state\ of\ reply\ posts]\ ([Sentiment\ state\ of\ intermediate\ self-reply]\ [Average\ sentiment\ state\ of\ reply\ posts])^*\ [Sentiment\ state\ of\ final\ self-reply]$  where average sentiment state of reply posts is obtained from the average sentiment probability defined in equation (1) using the threshold of sentiment state classifier described in section 3.1 (i.e., 0.5). More precisely, each sentiment state can take one of two values: positive or negative. Let  $b$ ,  $o$ ,  $r$ ,  $s$ , and  $f$  be atomic propositions. Let  $b$  denote the beginning of a thread,  $o$  denote that the initial post sentiment is positive,  $r$  denote that the average sentiment of reply posts elicited by the initial post is positive,  $s$  denote that the sentiment of an intermediate self-reply to the initial post is positive,  $f$  denote that the sentiment of the final self-reply is positive. A thread can be represented by a sentiment state sequence  $\mathbf{x} = x_0x_1 \cdots x_n$  where  $x_i \in \mathcal{X} = \{o, \neg o, r, \neg r, s, \neg s, f, \neg f\}$  where  $\neg$  denotes negation of the corresponding proposition.

### 3.3 Probabilistic Kripke Structure in CSN

We use a probabilistic Kripke structure [19] to represent and reason about probabilistic transitions between the sentiment states of posts in a CSN. Let  $\mathbf{x} = x_0x_1 \cdots x_n$  be a sequence of post sentiments where  $x_i \in \mathcal{X}$  and let  $X_i$  ( $0 \leq i \leq n$ ) be random variable corresponding to sequence element  $x_i$ . Markov Model (MM), which captures the dependencies between the neighboring sequence elements, is used to estimate the transition probabilities between sentiment states of posts that make up a thread. In  $k^{th}$  order Markov model, the sequence element follows the Markov property:  $X_i \perp\!\!\!\perp \{X_0, \dots, X_{i-k-1}\} \mid \{X_{i-k}, \dots, X_{i-1}\}$  (i.e.,  $X_i$  is conditionally independent of  $X_0, \dots, X_{i-k-1}$  given  $X_{i-k}, \dots, X_{i-1}$  for  $i = k, \dots, n$ ). Formally, the transition probabilities are estimated<sup>4</sup> over a set  $\mathcal{D} = \{\mathbf{x}_l\}_{l=1}^{|\mathcal{D}|}$  of sentiment sequences as follows.

$$\hat{p}(X_i|w) = \left[ \frac{1 + \sum_{l=1}^{|\mathcal{D}|} \# [w\sigma, \mathbf{x}_l]}{|\mathcal{X}| + \sum_{\sigma' \in \mathcal{X}} \sum_{l=1}^{|\mathcal{D}|} \# [w\sigma', \mathbf{x}_l]} \right]_{\sigma \in \mathcal{X}} \quad (2)$$

where  $\# [w\sigma, \mathbf{x}_l]$  represents the number of times the symbol  $\sigma$  “follows” the subsequence  $w$  (of length  $k$ ) in sequence  $\mathbf{x}_l$  and  $\hat{p}(X_i|w)$  is the estimate of the conditional

<sup>4</sup> Laplace correction is used for smoothing purposes

probability  $P(X_i = \sigma | w)$  of sequence element  $X_i$  that “follows” the subsequence  $w$ . We use the first-order MM to determine the transition probabilities for the CSN probabilistic Kripke structure.

## 4 Temporal Causality in Cancer Survivor Network

### 4.1 Prima Facie Cause

Figure 2 shows the Probabilistic Kripke structure ( $K$ ) that is constructed using the method described in previous section. The structure shows that from any state of the thread originator, i.e.,  $\{o, \neg o, s, \neg s\}$ , there is a probability greater than 74% that it will transit to the state  $r$ . This suggests that people who reply to thread originators have a

Table 2: Sentiment Dynamics in CSN

# $[o]$	# $[\neg o]$	# $[\neg o \rightarrow f]$	# $[o \rightarrow \neg f]$
12147	10707	7512	2948

high tendency to express positive sentiment regardless the sentiment of the thread originators at the beginning and in the middle of the thread. In other words, members of CSN try to offer positive social support to others who seek support. The thread originator starts with either initial positive or negative sentiment and after the interaction with other people in the CSN forum he/she might end up with either positive or negative sentiment. Table 2 shows that about 72% of thread originators with initial negative sentiment end up with positive sentiment, and only about 24% of thread originators with initial positive sentiment end up with negative sentiment at the end.

Table 3: Yearly Prima Facie Causes of Final Thread Originator’s Sentiment

Year	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010
$f$	$o$	$o, r$	$r, s$	$r, s$	$r, s$	$r, s$	$r, s$	$r, s$	$r$	$r, s$	$r, s$
$\neg f$	$\neg o$	$\neg o, \neg r, \neg s$	$\neg o, \neg r$	$\neg r, \neg s$	$\neg r, \neg s$	$\neg r, \neg s$	$\neg r, \neg s$	$\neg r, \neg s$	$\neg r$	$\neg r, \neg s$	$\neg r, \neg s$

Table 4: Communal Prima Facie Causes of Final Thread Originator’s Sentiment

Community	Breast Cancer	Colorectal Cancer
$f$	$r$	$r$
$\neg f$	$\neg r$	$\neg r, \neg s$

Our goal is to uncover the *prima facie* causes for final sentiment of the thread originators. Based on the definition of *prima facie* causes and the probabilistic Kripke structure  $K$ , we find that the set of *prima facie* causes of  $f$  and  $\neg f$  are  $\{r, s\}$  and  $\{\neg o, \neg r, \neg s\}$ , respectively.

CSN is comprised of users’ data over a period of 11 years. CSN can be divided into several sub-communities (e.g., Breast Cancer, Colorectal Cancer). To validate the above *prima facie* causes, we divided CSN data set into several subsets based on the year and the sub-community.

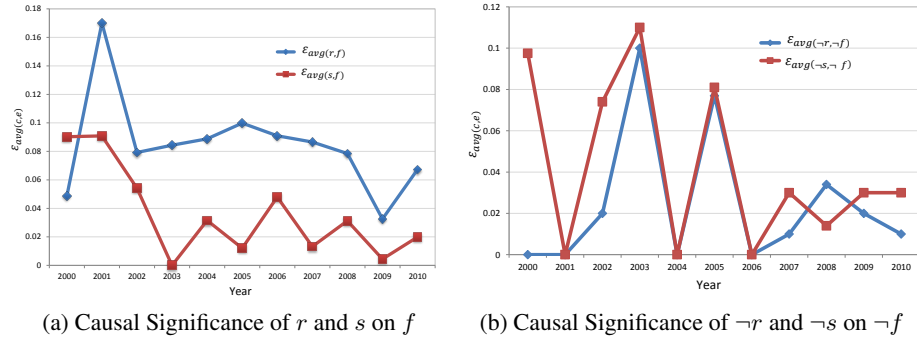


Fig. 3: Causal Significance According to Year.

Specifically, we group threads which were started in the same year (from 2000 to 2010) and we group threads that belong to either Breast Cancer or Colorectal Cancer community (during the period from 2005 to 2010). Surprisingly,  $r$  and  $\neg r$  are consistently the prima facie causes of  $f$  and  $\neg f$ , respectively in both yearly and sub-community data sets. Table 3 and 4 show the prima facie causes of  $f$  and  $\neg f$  in the yearly and sub-community analysis. *The results from the two tables indicate that the positive sentiment of the replies appears to causally influence the positive sentiment of the thread originator at the end of the thread; and conversely, the negative sentiment of the replies appear to causally influence the negative sentiment of the thread originators at the end of the thread.*

#### 4.2 Assessing the significance of causes

We proceed to evaluate the significance of the prima facie causes of  $f$  and  $\neg f$  using the method introduced in [21]. We calculate the significance of a cause  $c$  for an effect  $e$  as  $\varepsilon_{avg}(c, e) = \frac{\sum_{x \in X \setminus c} \varepsilon_x(c, e)}{|X \setminus c|}$ , where  $X$  is a set of prima facie causes of  $e$  and  $\varepsilon_x(c, e) = P(e|c \wedge x) - P(e|\neg c \wedge x)$  denotes the contribution of  $c$  to the change in probability of  $e$ . Table 5 shows causal significance between causes and effects from an aggregate of data from all the years.

From the table 5, we can see that causal significance  $\varepsilon_{avg}(r, f)$  is much higher than the causal significance  $\varepsilon_{avg}(s, f)$  and whereas  $\varepsilon_{avg}(\neg o, \neg f)$ ,  $\varepsilon_{avg}(\neg r, \neg f)$ , and  $\varepsilon_{avg}(\neg s, \neg f)$  are not much different from each other.

Table 5: Causal Significance

$\varepsilon_{avg}(r, f)$	$\varepsilon_{avg}(s, f)$	$\varepsilon_{avg}(\neg o, \neg f)$	$\varepsilon_{avg}(\neg r, \neg f)$	$\varepsilon_{avg}(\neg s, \neg f)$
0.054	0.01	0.05	0.04	0.039

In a similar fashion, we also examined the causal significance on data for specific years and sub-communities. Figure 3 shows the results of this analysis. Figure 3a shows that the causal significance  $\varepsilon_{avg}(r, f)$  is significantly greater (pair t-test,  $p < 0.01$ ) than the causal significance  $\varepsilon_{avg}(s, f)$ . However, figure 3b shows that the  $\varepsilon_{avg}(\neg r, \neg f)$  and



$\varepsilon_{avg}(\neg s, \neg f)$  are not significantly different (figure 3b does not include the significance of  $\neg o$  since it is not found to be a cause of  $\neg f$  in most of the years (except during the first three years which account for less than 4% of the total number threads)). Our analysis of the data from the sub-communities yields a similar finding (i.e.,  $\varepsilon_{avg}(r, f)$  is significantly greater than  $\varepsilon_{avg}(s, f)$  and  $\varepsilon_{avg}(\neg r, \neg f)$  and  $\varepsilon_{avg}(\neg s, \neg f)$  are not significantly different from each other).

Based on the results summarized in table 5 and figure 3, we can conclude that  $r$  causally influences  $f$  and  $\{\neg r, \neg s\}$  causally affect  $\neg f$ . In other words, our key finding is that the positive sentiment of a reply causes the negative to positive change in the thread originator’s sentiment, at least among 72% of the thread originators with initial negative sentiment. We also see that negative sentiment from a replier causes the thread originator to be left with a negative sentiment. Hence, we conclude that the sentiment of the replies drives the sentiment dynamics of the thread originator.

## 5 Summary and Future Work

In this work, we have introduced a framework to uncover the temporal causality of sentiment dynamics of the thread originator in the American Cancer Society’s Cancer Survivor Network. To the best of our knowledge, this study is the first to uncover the factors that causally drive the sentiment dynamics in an OHC. We developed a sentiment classifier using machine learning on a training set of posts with manually labeled for their sentiment (positive versus negative). We constructed a Probabilistic Computation Tree Logic representation and a corresponding probabilistic Kripke structure to represent and reason about the transitions between sentiments of posts in a thread over time. We analyzed the Kripke structure to identify the prima facie causes of sentiment change on the part of the thread originators in the CSN forum and their significance. *Our main finding is that the positive sentiment of replies appears to causally influence the positive sentiment of the thread originator at the end of the thread; and conversely, the negative sentiment of the replies appears to causally influence the negative sentiment of the thread originators at the end of the thread.* Some promising directions for future research include: (i) exploring the causal effects of the topic being discussed on the sentiment dynamics; (ii) exploring the causal effects of (explicit as well as implicit) social relations among OHC participants on sentiment dynamics.

**Acknowledgements:** This research is supported by a Collaborative Agreement with American Cancer Society, which made the data of CSN available for this research. The work of Bui was supported in part by a research assistantship funded by the Edward Frymoyer Endowed Chair in Information Sciences and Technology held by Vasant Honavar. This work has benefited from discussions with Kenneth Portier, Greta Greer (both at American Cancer Society), Prasenjit Mitra (Qatar Computing Research Institute), Cornelia Caragea (University of North Texas), Kang Zhao (University of Iowa), David Reitter, and other current and former members of Cancer Informatics group at Penn State University.

## References

1. Ferlay, J., et al.: Cancer incidence and mortality worldwide: sources, methods and major patterns in globocan 2012. *International Journal of Cancer* (2014)
2. American Cancer Society: American cancer society cancer facts and figures 2014. <http://www.cancer.org/research/cancerfactsstatistics/cancerfactsfigures2014/> (2014)
3. Fox, S., Duggan, M.: Health online 2013. Pew Research Internet Report (2013)
4. Dunkel-Schetter, C.: Social support and cancer: Findings based on patient interviews and their implications. *Journal of Social Issues* **40** (1984) 77–98
5. Preece, J.: Empathic communities: balancing emotional and factual communication. *Interacting with Computers* **12** (1999) 63–77
6. Rodgers, S., Chen, Q.: Internet community group participation: Psychosocial benefits for women with breast cancer. *Journal of Computer-Mediated Communication* **10**(4) (July 2005)
7. Beaudoin, C.E., Tao, C.C.: Modeling the impact of online cancer resources on supporters of cancer patients. *New Media and Society* **10**(2) (2008) 321–344
8. Maloney-Krichmar, D., Preece, J.: A multilevel analysis of sociability, usability, and community dynamics in an online health community. *ACM Trans. Comput.-Hum. Interact.* **12**(2) (June 2005) 201–232
9. Portier, K., Greer, G.E., Rokach, L., Ofek, N., Wang, Y., Biyani, P., Yu, M., Banerjee, S., Zhao, K., Mitra, P., Yen, J.: Understanding topics and sentiment in an online cancer survivor community. *Journal of the National Cancer Institute (JNCI) Monographs* (2013) 195–198
10. Qiu, B., Zhao, K., Mitra, P., Wu, D., Caragea, C., Yen, J., Greer, G.E., Portier, K.: Get online support, feel better - sentiment analysis and dynamics in an online cancer survivor community. In: *SocialCom/PASSAT*. (2011) 274–281
11. Huh, J., Yetisgen-Yildiz, M., Pratt, W.: Text classification for assisting moderators in online health communities. *J. of Biomedical Informatics* **46**(6) (December 2013) 998–1005
12. Biyani, P., Caragea, C., Mitra, P., Zhou, C., Yen, J., Greer, G.E., Portier, K.: Co-training over domain-independent and domain-dependent features for sentiment analysis of an online cancer support community. *ASONAM '13*, New York, NY, USA, ACM (2013) 413–417
13. Wang, X., Zhao, K., Street, N.: Social support and user engagement in online health communities. In: *Smart Health*. Volume 8549 of *Lecture Notes in Computer Science*. Springer International Publishing (2014) 97–110
14. Kleinberg, S.: *Causality, Probability, and Time*. Cambridge University Press (2013)
15. Kleinberg, S., Hripacsak, G.: A review of causal inference for biomedical informatics. *Journal of Biomedical Informatics* (2011) 1102–1112
16. Pearl, J.: *Causality: Models, Reasoning, and Inference*. Cambridge University Press (2000)
17. Suppes, P.: *A Probabilistic Theory of Causality*. Noth-Holland, Amsterdam (1970)
18. Prior, A.: *Past, Present, and Future*. Clarendon Press, Oxford (1967)
19. Clarke, Jr., E.M., Grumberg, O., Peled, D.A.: *Model Checking*. MIT Press, Cambridge, MA, USA (1999)
20. Hansson, H., Jonsson, B.: A logic for reasoning about time and reliability. *Formal Aspects of Computing* **6** (1994) 102–111
21. Kleinberg, S., Mishra, B.: The temporal logic of causal structures. In *Proceeding of the UAI '09*, Arlington, Virginia, United States, AUAI Press (2009) 303–312
22. Thelwall, M., Buckley, K., Paltoglou, G., Cai, D., Kappas, A.: Sentiment in short strength detection informal text. *J. Am. Soc. Inf. Sci. Technol.* **61**(12) (December 2010) 2544–2558
23. Hu, M., Liu, B.: Mining and summarizing customer reviews. In *Proceeding of KDD '04*, New York, NY, USA, ACM (2004) 168–177