

Activating the "Breakfast Club": Modeling Influence Spread in Natural-World Social Networks

Lily Hu¹, Bryan Wilder², Amulya Yadav², Eric Rice², Milind Tambe²

¹Harvard University ²University of Southern California

lilyhu@g.harvard.edu, {[bwilder](mailto:bwilder@usc.edu),[amulyaya](mailto:amulyaya@usc.edu),[eric](mailto:eric@usc.edu),[tambe](mailto:tambe@usc.edu)}@usc.edu

ABSTRACT

While reigning models of diffusion have privileged the structure of a given social network as the key to informational exchange, real human interactions do not appear to take place on a single graph of connections. Using data collected from a pilot study of the spread of HIV awareness in social networks of homeless youth, we show that health information *did not* diffuse in the field according to the processes outlined by dominant models. Since physical network diffusion scenarios often diverge from their more well-studied counterparts on digital networks, we propose an alternative Activation Jump Model (AJM) that describes information diffusion on physical networks from a *multi-agent team* perspective. Our model exhibits two main differentiating features from leading cascade and threshold models of influence spread: 1) The structural composition of a seed set team impacts each individual node's influencing behavior, and 2) an influencing node may spread information to non-neighbors. We show that the AJM significantly outperforms existing models in its fit to the observed node-level influence data on the youth networks. We then prove theoretical results, showing that the AJM exhibits many well-behaved properties shared by dominant models. Our results suggest that the AJM presents a flexible and more accurate model of network diffusion that may better inform influence maximization in the field.

KEYWORDS

Social networks; team formation; innovative applications

1 INTRODUCTION

To say that we reside in networks may yet be an understatement. For as long as human beings have been social beings, we have been embedded in connectivity. Our mesh of interpersonal ties contains both relational and informational content. Networks reveal not only *who* we may know but also *what* we may know and *when* we may come to know it. Research in influence maximization lies at this crossroads of the *who*, *what*, and *when* of information diffusion. In their seminal paper, Kempe, Kleinberg, and Tardos [23] formalized the problem by imagining influencing agents as seed nodes in a network initialized to propagate information first to neighbors and eventually throughout the network as influence spreads. Along with theoretical work in the field, the wide availability of large-scale digital data has positioned internet networks—web link traces, email communication, social media platforms—at the center of the discussion of the

influence maximization problem. Even so, empirical work that compares predictions by the dominant Independent Cascade (ICM) and Linear Threshold (LTM) models with realized diffusion at the *node level* is limited. Prevailing methodologies for estimating diffusion model parameters often achieve low accuracy in replicating observed behavior even when applied to well-defined online networks with temporal information flow data [17, 38].

Although the ICM and LTM were originally formulated to describe social influence in natural environments (environments where agents physically interact, as opposed to purely online ones) [18, 34], there is a dearth of high-quality data that support the theories on networks in physical settings. Moreover, many of the starting premises of the leading information diffusion models are difficult to generalize to physical settings, presenting challenges that exacerbate the data deficiencies. First, both the ICM and LTM assume that the topology of social ties is identical with the mesh of connective channels through which information spreads. While on many social media platforms a user's social network delimits her space of communication, in the natural world, an individual's total space of social navigation dwarfs the space of those she calls her "friends," and there exists a multiplicity of information avenues within the network that do not coincide with one's social ties [35]¹. Further, previous work has shown that models with strong assumptions about a particular constructed graph topology are more prone to error and inaccuracy in their predictions of information spread [7, 24].

In this paper, we analyze information diffusion data from the first empirical study of influence maximization in the *physical world* [39], which tracked 173 individuals across 3 distinct networks over a multi-year time period. We delved into this significant corpus of natural world network data to investigate influence spread at an *individual node level* rather than a network-wide volume level. We found that information *did not* diffuse from seed nodes to the greater network according to processes suggested by the Independent Cascade or Linear Threshold models. Most strikingly, we found that across all three networks in the study, *50% of informed nodes lacked any path to a seed node* and moreover, a node's degree of connectivity—both generally and specifically to seed nodes—exhibited *no correlation with likelihood of becoming informed*. These results directly contradict predictions put forth by the ICM and LTM and call into question the suitability of these leading models of diffusion for approximating information spread on physical networks.

These negative results against the ICM and LTM can be generalized to apply to other graph-based models of diffusion: The high proportion of informed yet isolated nodes cannot be explained by models that rely on edge-based propagation as the sole avenue for

Proc. of the 17th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2018), M. Dastani, G. Sukthankar, E. Andre, S. Koenig (eds.), July 2018, Stockholm, Sweden

© 2018 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.
<https://doi.org/doi>

¹Prior work has suggested that nodes may receive information via "external influence" originating from outside the network [31]; here we suggest that nodes may still be influenced from *within* the network via a non-neighboring node.

influence spread [9]. On the other hand, more flexible approaches such as Hawkes contagion processes [40], which allow for influence to spread over greater social distances, cannot explain the differing spread outcomes across the networks. In general, appealing only to non-graph features of the diffusion scenarios cannot reconcile the divergences in observed information diffusion in the three different Yadav et al. studies. In response to the shortcomings of existing models, we develop a new model of *targeted peer-to-peer information spread on natural networks* that does not rely on strong tie assumptions and instead incorporates an understanding of influencing as a “team” behavior.

Our model features two distinguishing characteristics that are aligned with real-world information diffusion in physical spaces: 1) Nodes exchange information beyond their immediate social ties, and 2) Seed nodes act as a *multi-agent team* to spread information, where their overall influencing efficacy is a function of both individual and team attributes. In the proposed *Activation Jump Model* (AJM), team-based influence spread in a network is driven by activating the “Breakfast Club,”² where individuals from different social contexts band together for a common cause and form a united team for information diffusion. These features confer a flexibility to our forecasts of information flow and allow our model to *achieve a 60% to 110% improvement over the best ICM and LTM predictions in its predictions of which nodes will be influenced*. Although there is a long line of multi-agent systems work on modeling social influence [15, 20, 27], none has considered these aspects of natural-world information diffusion.

We also point to a methodological pitfall of research in influence maximization that focuses solely on achieving a particular level of information diffusion within a network. Namely, matching magnitude of influence spread under simulations to observed influence spread is insufficient evidence for determining the underlying diffusion process. We show that optimal seeding under one model achieves *near-optimal* (> 90%) influence spread under another diffusion process on three natural-world networks. In fact, *any* magnitude of influence spread can be explained by varying ICM and LTM parameters, pointing to a fundamental ambiguity in identifying the true diffusion process based on this metric alone.

Even when seeding strategies achieve high levels of influence spread, leading models’ failures to predict node-level influence can limit their applicability for real-world uses. In domains of sustainability, network interventions can generate knowledge as well as promote behavioral changes within a community. These programs typically identify individuals and groups that may especially benefit from the intervention. For example, school network-based suicide prevention programs aim to increase general awareness about signs of suicidal behavior but especially seek to reach high-risk adolescents and their social circles [22]. Similarly, peer-led HIV prevention programs akin to the fieldwork by Yadav et al. [39] hope to reach a diverse set of individuals but especially those who participate in risky behaviors that increase their likelihood of HIV infection [5]. As such, most social interventions have the dual purpose of maximizing influence coverage while also targeting particular vulnerable individuals. Furthermore, such social programs are typically deployed in

under-resourced communities within which the reach of digital networks may be more limited and information exchange predominantly occurs on physical social networks. With its superior performance in predicting node-level influence in these arenas, the AJM may serve as a more desirable framework to guide these intervention strategies in the field.

In the next section, we present results from our study of influence maximization on physical networks and show that the observed patterns of informational exchange do not accord with predictions by current leading models. In Section 3, we propose an alternative model of information diffusion that more closely approximates true influencing behavior in the field and exhibits improved prediction performance of node-level information spread. We then prove theoretical results relating to the influence maximization problem under the new model in Section 4. The paper concludes with a Discussion section that looks toward potential applications for our model in socially-oriented domains.

2 NETWORK INTERVENTION DATA ANALYSIS

2.1 Pilot Study Procedure

Yadav et al.’s long-standing collaboration with homeless youth service providers in a large urban area sought to improve peer-led health interventions by leveraging research in influence maximization ([2017]). To this end, they conducted a series of head-to-head comparison studies of various seeding strategies to select cohorts of *Peer Leaders* among the youth that would be trained for the task of HIV awareness diffusion in their communities. Three studies took place on three distinct social networks of homeless youth. Each study recruited youth and gathered social network data using online contacts, field observations, and surveys. A different seeding strategy was then deployed on each of the generated networks: In two of the pilot studies, Peer Leaders were chosen via two algorithmic agents for influence maximization, HEALER and DOSIM, which were designed to optimize network-based intervention strategies for health providers. The third network was seeded via degree centrality (DC), the most commonly-used heuristic in network interventions [37], which simply selected the most popular youth, those with the greatest number of social connections, to become Peer Leaders.

Each network’s Peer Leaders underwent an intensive training course led by pilot study staff that served to both instruct the youth in spreading information about HIV to their peers as well as bind the members together in their shared roles as health ambassadors. After Peer Leaders were sent out into the field, youth were asked in 1-month and 3-month follow-up surveys about whether they had received information about HIV from a Peer Leader. These responses revealed the extent to which information had spread from seed nodes to the greater network. The post-intervention results revealed that the HEALER and DOSIM seeding strategies resulted in greater informational spread compared to DC, with ~74% and ~72% respectively of non-Peer Leaders reporting having received information about HIV in the 3-month survey compared to ~35% in the control study. Since both HEALER and DOSIM solved the influence maximization problem by assuming a model of information spread based on a generalization of the Independent Cascade, the success initially seemed

²<http://www.imdb.com/title/tt0088847/>

Table 1: Fraction of nodes in each network that were informed, sorted by connectivity status to Peer Leaders. Denominator gives total number of nodes of that type; numerator gives number of those that are informed. Direct nodes have an edge to a PL; indirect nodes are connected via intervening neighbor(s); isolated nodes lack a path to any PL.

Network	n	Direct	Indirect	Isolated	Proportion Informed
HEALER	34	15/21	4/7	6/6	25/34
DOSIM	25	5/6	5/10	8/9	18/25
DC	26	5/12	1/4	3/10	9/26

to validate the model as an accurate approximation of information spread in the physical world.

2.2 Node-level analysis of information diffusion

However, the empirically observed *node-level* patterns of information spread in the three networks wildly diverged from Independent Cascade and Linear Threshold predictions. Table 1 gives an overview of the connectivity of nodes that reported receiving information about HIV from a Peer Leader in the 3-month follow-up survey. In each of the three networks, nodes lacking a path to any seed Peer Leaders—denoted as “isolated” in the table—represented a high proportion of all nodes that were informed. These “isolated” nodes either occupied a distinct connected component in the graph separate from all Peer Leaders or had no edges entirely. Whereas under the Independent Cascade and Linear Threshold models, these types of nodes would have a 0 probability of being informed, notably, in both the HEALER and DOSIM interventions, isolated youth were informed at a rate *higher* than even those youth who were directly connected to one or more Peer Leaders, with 100% (6/6 in HEALER) and 89% (8/9 in DOSIM) informed compared to ~71% (15/21) and ~83% (5/6). In the DC network, the effect is less pronounced, though isolated nodes were still informed at a rate comparable to the general non-Peer Leader population (30% compared to ~35%). Nonetheless within the context of the ICM and LTM, such nodes have a 0 probability of receiving information. Thus, these results immediately challenge the claim that existing ties are the dominant avenues of informational exchange and also call into question the premise that information radiates out from seed nodes first to neighbors and then to the rest of the network.

In order to more finely assess the effect that a node’s connectivity had on its likelihood of receiving HIV information, we calculated Pearson correlation coefficients between two degree measures and a node’s final information status. Our results in Table 2 show that all such correlations are not significantly different from no correlation, thus indicating that connectivity has *no bearing* on likelihood of being influenced. This stands in contrast to prevailing models, in which a node’s edges represent its “opportunities” to receive information, and thus both Peer Leader degree—the number of ties a node has to Peer Leaders—and total degree should be strictly positively correlated with becoming informed.

The lack of positive correlation between either total degree or number of connections with Peer Leaders and influence status is even more dissonant with edge-based models of propagation when

Table 2: Pearson correlation coefficients r between nodes’ influence statuses and their Peer Leader (PL) and total degree. PL degree counts edges to PLs; total degree counts edges to all nodes. Positive (negative) r implies a positive (negative) relationship between reporting hearing HIV information and degree. 95% confidence intervals are also given.

Network	n	PL Degree	Total Degree
HEALER	34	-0.0685 (-0.397, 0.276)	-0.1867 (-0.494, 0.162)
DOSIM	25	0.1418 (-0.268, 0.508)	-0.0290 (-0.419, 0.370)
DC	26	0.1547 (-0.247, 0.511)	0.1304 (-0.271, 0.493)

considered alongside the high levels of influence spread achieved in the studies. In the HEALER network, information was successfully transmitted to ~74% of all non-Peer Leaders, corresponding to a most likely propagation probability of $p \approx 0.84$. Such a high propagation probability further suggests that Peer Leader-neighboring nodes should be even more heavily favored to receive information, with simulations predicting that *nearly all* (~99%) would become informed, whereas in reality, only ~71% of these nodes received information. Simulations on the HEALER network with this p value produce correlations of 0.489 and 0.598 between degree and likelihood of being informed (PL and total respectively), indicating a moderate to strong positive relationship compared to the actual values of -0.0685 and -0.1867 , which indicate negative to no relationship between degree and influence status. In the DOSIM study, the graph topology itself, with 5 connected components in addition to 7 nodes of degree 0, restricts information spread under the ICM and LTM to maximally reach 68% of all non-Peer Leaders. Even under perfect information propagation, as long as nodes are only able to influence neighbors, simulations under-predict the observed information spread.

Discussion of alternate hypotheses: We now address alternate explanations that may be suggested to explain why these results greatly diverge from predictions made by graph-based models like the ICM and LTM. We draw upon a large body of existing work in sociology and network theory to argue that such alternate hypotheses are highly implausible or at the very least insufficient to explain the observed pattern of influence spread.

We start with the prominent presence of isolated nodes, as their high rate of conversion is one of greatest points of dissonance with predictions put forth by the ICM and LTM. One counter-hypothesis may suggest that these nodes were not truly isolated but in fact had social ties that had gone unreported and as a result, were not captured in the graph collection process. However, we note that the existence of truly isolated nodes is consistent with previous social work with adolescent populations. Networks composed of a large connected component along with several dyads, triads, and a large number of isolates have been confirmed in work on high-risk adolescents in the Bronx, other work with homeless youth communities, and friendship networks among high school students [12, 33, 37]. In general, when there is a loose boundary to a community such as a school, homeless drop-in center, or neighborhood, adolescents form social networks with a fairly consistent structure akin to the

graphs collected in Yadav et al. [39]. Social isolates are therefore not aberrations of these types of networks, rather they are *features*. We conclude that dismissing *prima facie* the existence of truly isolated nodes and assuming that such nodes are “mistakenly” isolated, is not viable within the larger context and lineage of research in these areas, and thus (at least many of) these nodes are truly isolated and had no preexisting social connections.

Further, in the follow-up survey, such youth definitively identified having had a conversation with a Peer Leader. The survey stated, “For the next several questions we are going to ask you about conversations you may have had about HIV or AIDS in the past month with a Peer Advocate in the Have You Heard program.”, followed by a list of the Peer Leaders’ names. Youth were then asked to rate their agreement with the statement “The conversation was awkward” and were given the option of selecting the response “I didn’t have a conversation with a Have You Heard Peer Advocate about HIV or AIDS.” if no such Peer Leader conversation had occurred. *The isolated nodes in Table 1 responded in the affirmative to the question, answering how they felt about the conversation (rather than indicating that no such conversation occurred)*. Thus, these responses provide direct evidence that social influence occurred outside of preexisting network connections.

We now turn to the finding that there is no significant correlation between a node’s degree of connection with Peer Leaders and their likelihood of becoming informed. An alternative hypothesis for such a finding could be that the constructed network was missing a crucial set of connecting edges along which information actually diffused. But in fact, research on social network data collection indicate that collected graphs are actually more likely to be systematically skewed *in favor* of the ICM and LTM. The standard graph collection techniques that were used in the field study are biased toward persons reporting edges where there is greater emotional strength or more frequent interactions [4, 29], and these are precisely the strong, frequently accessed ties which would be most likely to propagate information under prevailing models. Hence, our dataset should be biased to *include* edges used to spread information and omit weaker ties that are less likely to carry information. While it is always certainly possible to collect more detailed data in the field, it is highly unlikely that all inaccuracies in these current methods were aligned *exactly* to invalidate the ICM and LTM (and occurred similarly in three separate networks). Thus even on a constructed network that should have been biased *towards* the ICM and LTM, we nevertheless see no evidence that diffusion is explained solely by existing ties. There is therefore strong evidence that an alternate mechanism of information spread underlies the observed outcomes.

These results are also consistent with prominent social network theories such as sociologists Mark Granovetter’s and Ronald Burt’s theories of weak ties [19] and structural holes [6]. Both are based upon empirical findings that suggest that social influence often occurs in counter-intuitive and non-linear ways. In each of these theories, influence occurs such that more distant or indirect social ties can have greater importance than strong direct ties, especially when novel information is being disseminated within a network, which stands in direct opposition to leading graph-based computational models of social influence today. We see our new model, proposed in the following section, as a formal mathematical treatment of similar social dynamics as theorized by Granovetter and Burt.

Taken together, the multiple aforementioned contradictions with prevailing models indicate that the empirical results cannot be dismissed as simply anomalous. Given the unique challenges and complexities of information diffusion on physical networks, we suggest that the data’s divergence from predictions by models that have been largely validated only on digital networks is one step in uncovering and understanding a qualitatively different influence process. We thus conclude that there is no evidence that a cascade or threshold-like process of information diffusion produced the observed data and move toward developing a new model of influence dynamics on real-world physical networks.

3 PROPOSED MODEL

In this section, we introduce a new model of information spread for this class of peer-to-peer diffusion phenomena.

3.1 Activation Jump Model

Beginning with the premise that instances of social exchange are not limited to nodes that share a tie, our model of diffusion does not constrain information flow to the edges within a network. In the *Activation Jump Model* (AJM), influencing agents may leave their immediate social neighborhood to contact and propagate information to other nodes. This action of contacting nodes beyond one’s first-order ties is signified as a “jump.” We recognize the heterogeneity of active nodes’ social dispositions by differentially modeling each influencer’s jump behavior. A seed node’s jump activity has two main components: 1) *activation level*, a measure of *how many* other nodes it will attempt to influence throughout the course of the diffusion period, and 2) *landing distribution*, a probability distribution that expresses to *which* inactive nodes it will jump. Together, these two features describe *how often* and *to whom* an influencing agent contacts as she navigates the network to spread information.

Thus the AJM comprises two stages: First, each seed node determines its activation level, giving the number of other nodes to which it will jump. Second, the seed set is deployed in the network, and the social influence process unfolds in time. When a given seed jumps at time t , it selects from its landing distribution a target node uninformed at time t , modeling the process by which seed nodes seek nodes to inform. Influence is then successfully propagated to an uninformed node with probability p .

The AJM takes a multi-agent systems approach to the influence maximization problem by constructing a model of node activation that is a function of both individual and “team” attributes. In contrast to prior models, the seed set is not a collection of independent influencers, rather nodes exhibit behavioral dependencies wherein group dynamics either contribute to or detract from aggregate activation levels.

3.2 Model Formalization

While throughout this paper and in all our results, we use a form of the AJM tied to the graph’s structural properties, we first discuss the general form of the model to show that it can accommodate a broad range of properties and then discuss our specific form. We return to the generalization of the model in the Discussion.

Consider a team of seed nodes, S , tasked with information diffusion on a network $G = (V, E)$. Each seed node draws its activation

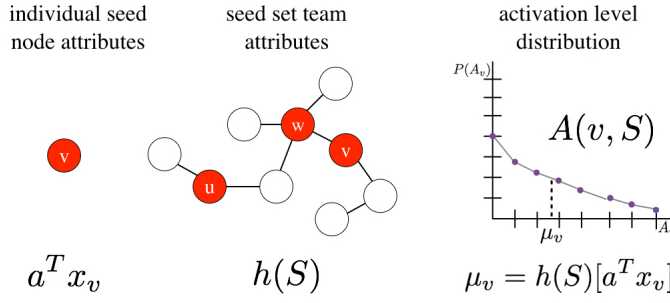


Figure 1: In the Activation Jump Model, a seed node v is associated with an activation level distribution that is a function of individual node as well as seed set S attributes.

level, giving the number of social jumps, or influencing attempts, it will make, from a distribution that is a function of both the node's *individual* attributes as well as the seed set's *team* attributes. Formally, let Δ^Z be the set of distributions over integers $Z \geq 0$. Each node $v \in V$ is associated with a function $f_v : 2^{|V|} \rightarrow \Delta^Z$ that maps the set of seed nodes to a distribution $A(v, S)$ over discrete activation levels. $A(v, S)$ is a parameterized distribution (e.g., geometric) with mean $\mu_v = h(S)[a^T x_v]$ where x_v is the node's attribute vector with coefficients a^T . Together $a^T x_v$ represents the particular node's maximum activation level, which is modulated by the team activation level term given by $h(S) \leq 1$, a function of the structural positions of nodes in S that captures discomplementarities among team members. Figure 1 illustrates this dual—individual and team—composition of a seed node's activation level distribution.

Each node v is also associated with a *landing distribution*, $L_{v,T}$, giving the probabilities with which v jumps to a set of potential target nodes T . The landing probability is a function of the attributes of the influencing seed v and particular targeted node u . Based on these qualities, the node pair is assigned a score $\phi(v, u) \geq 0$, and $L_{v,T}(u) = \Gamma \phi(v, u)$ where Γ is a normalization factor such that $\sum_{u \in T} L_{v,T}(u) = 1$.

We now instantiate the AJM in a specific form that features the concept of “structural diversity,” which highlights groups with members who participate in multiple distinct social contexts. The team thus acts to *unite* otherwise disparate nodes, producing the “Breakfast Club” effect, which has been shown to be a key determinant of diffusion in networks [36]. Thus we formulate the function

$$h(S) = 1 - \frac{1}{k} \sum_{(u,v) \in E} \mathbb{1}[u, v \in S] \quad (1)$$

where k is the total number of seed nodes to be selected. h illustrates the negative effect of social homogeneity in the form of between-seed-node edges on a team's effectiveness. Each pair of connected seed nodes entails a loss of $\frac{1}{k}$ of the team's effectiveness. Barring negative influence, $h'(S) = \max(h(S), 0)$ without loss of generality. To reflect the correlation between degree and propensity towards sociality and thus activation, we parsimoniously set $a^T x_v = \text{deg}(v)$. That is, high-degree nodes will be more active in spreading influence (so long as they are part of a diverse team).

The landing distribution score for seed v and target u is given by $\phi(v, u) = \frac{1}{d(v, u)}$, where $d(v, u)$ is the path-length distance between

the two nodes. When $d(v, u) = \infty$ (there is no path from v to u), we set $\phi(v, u) = \epsilon$, where $\epsilon > 0$ is a small constant. As in the ICM, the propagation probability p is able to be varied (we discuss this further in Section 3.4).

Information diffusion thus occurs in two stages. First, during the *activation stage* each node in the seed set $v \in S$ is initialized by drawing an activation level A_v from its distribution $A(v, S)$. Then, the *jump stage* unfolds over the time interval $[0, 1]$. Each seed node v draws a series of jump times $t_1^v \dots t_{A_v}^v$ from the uniform distribution over $[0, 1]$. At each jump time t_i^v , v jumps to an uninfluenced target node u drawn from $L_{v,T}$ where T the set of uninfluenced nodes at t_i^v . Finally, u is successfully influenced with probability p .

3.3 Model Discussion

The Activation Jump Model's incorporation of seed set team dynamics follows a line of multi-agent systems research, which demonstrates the importance of careful team formation when agents must collaborate to achieve a goal [1, 13, 14, 26]. In particular, previous work has focused on the value of forming a *diverse* team [2, 21, 28]. In the AJM, we operationalize this concept by using the group effectiveness function h to model network structural diversity by penalizing seed sets with many within-team edges. Thus $h \in [0, 1]$ is decreasing in the level of connectivity among seed nodes, and influencing nodes are more active when they occupy distinct neighborhoods of the network rather than when the team is socially homogeneous.

It is important to note that under this model setup, a node's marginal effect on the aggregate activation level of a seed set is not guaranteed to be positive. There may exist a node $w \in V$ such that $\sum_{v \in S} h(S)[a^T x_v] > \sum_{v \in S \cup w} h(S \cup w)[a^T x_v]$, with the effect that the influence function $f(\cdot)$, giving the expected number of influenced nodes, is non-monotone. Although this is a significant departure from the ICM and LTM, we argue that non-monotonicity is a realistic feature of team-based influence spread. A new seed node may interfere with team dynamics, resulting in a deleterious effect that outweighs its positive individual contribution to the team. This balance between the quantity and quality of members in a seed set is an important consideration in team formation in the real world. As a result, the influence maximization problem under the AJM requires examination of not only a node's individual attributes but also its effect on the group composition of nodes already in the seed set.

We remark that the AJM is intended specifically to model information spread on physical, mid-scale networks. It is not appropriate for modeling “passive” diffusion processes such as the network spread dynamics of a disease, which may be transmitted by a non-seed node. Moreover, the AJM is a *progressive* model such that once a node becomes influenced, it cannot revert to being in an uninfluenced state. Hence, the model is not suited to situations where nodes may repeatedly change their mind depending on social circumstances. Nevertheless, this leaves a great space of possible information spread scenarios that may be captured by the AJM. The multi-agent aspect of the model makes it particularly amenable to modeling targeted campaigns run by a team of influencers, which are common whenever particular members with special knowledge are charged with

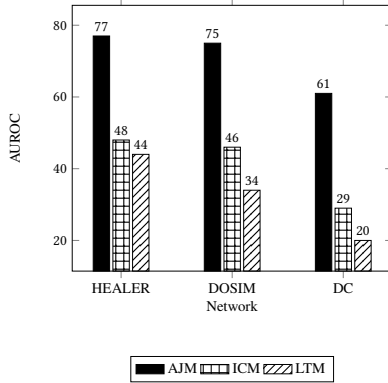


Figure 2: Comparison of model classification performances

disseminating social awareness within their communities. Such instances are common in the domains of public health, community organizing, and even team-based marketing.

3.4 Model Validation on Post-Intervention Data

We evaluate the performance of the Activation Jump Model by comparing its predictions to Yadav et al.’s dataset of HIV awareness spread on three distinct social networks of homeless youth. Standard experiments of diffusion models compare the magnitude of total influence spread under simulations to that observed empirically in order to assess model accuracy. Here, we perform a finer-grained analysis by evaluating and comparing AJM, ICM, and LTM predictions of node-level influence. We treat each model as a binary classifier that outputs the predicted probability of each node becoming influenced. Each model is then evaluated according to its AUROC, a standard measure of classification accuracy.

Parameter settings: Physical networks present challenges in data collection that limit the ability to view multiple cascades, rendering standard methods of inferring diffusion parameters inoperable. We work under this constraint by fitting the ICM directly to the test data by running simulations under the propagation value that gives its best classification performance and then forcing the AJM to also work under this propagation value. Thus any experimental bias in this setup favors the ICM.

For the AJM, the only parameter we set is the small constant $\phi(v, u) = 0.1$ for the landing distribution score when $d(v, u) = \infty$. By contrast, we present the strongest possible version of the ICM for each network, fitting it directly to the test data by selecting the propagation value p that maximized the ICM’s AUROC value. We then used this same probability for the AJM. By forcing the AJM to operate under the ICM’s optimal parameterization, we ensure that our experiments truly test the AJM’s better suitability for modeling the data, rather than a better ability to “memorize” the data.

Assessing classification accuracy: Using selected Peer Leaders in the field experiments as seed nodes, we generated diffusion instances according to the Activation Jump, Independent Cascade, and Linear Threshold models, tracing out Receiver Operating Characteristic (ROC) curves for each set of simulations. This evaluation methodology has been used in previous node-level analyses of information diffusion models [16, 38] and has been recognized as superior to

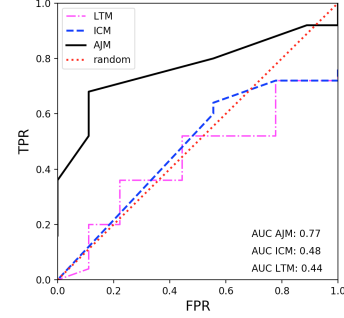


Figure 3: Models’ ROC curves on HEALER network

Precision-Recall curves for the binary classification task [32]. ROC curves plot a classifier’s True Positive Rate (TPR) against its False Positive Rate (FPR) with each point on the curve corresponding to a predictive threshold such that all nodes with a probability of being informed above (below) the threshold are classified as influenced (not influenced). We used the area under the ROC curve (AUROC) to evaluate classification performance [10] where an AUROC of 1 corresponds to a perfect classifier.

Results: Each model’s AUROC values for the three networks are shown in Figure 2; the ROC curves for all three models’ predictions on the HEALER network are shown in Figure 3. The AJM outperforms the ICM and LTM across all networks, with the model achieving accuracies (measured via AUROC) of 77% and 75% on the HEALER and DOSIM networks respectively, while the best ICM and LTM issue predictions that, on average, perform worse than a random classifier (20-48%). For DC, one possible explanation for all three models’ lower AUROCs is the overall poor permeation of influence throughout the network, since low base-rates cause the measure to be sensitive to small classification changes. Even so, the AJM is far from a trivial classifier, with an AUROC of 0.61 compared to the ICM and LTM values of 0.29 and 0.20 respectively.

4 THE INFLUENCE MAXIMIZATION PROBLEM UNDER THE ACTIVATION JUMP MODEL

We now consider the influencing activity of a coordinated multi-agent team under the AJM. Since seed agents do not target nodes that have already been informed, the influence function is captured by the total number of expected jumps, given by $f(S) = h(S) \sum_{v \in S} a^T x_v$, where $h(S)$ follows the form in Equation 1. We show that under natural conditions, f is a (potentially nonmonotone) submodular function.

LEMMA 4.1. *h is monotone-decreasing and submodular.*

PROOF. h is monotonically decreasing by inspection. For submodularity, consider the marginal impact of adding a given node v to an existing seed set S :

$$h(v|S) = -\frac{1}{k} \sum_{(u,w) \in E} \mathbb{1}[\{u,w\} \not\subseteq S, \{u,w\} \subseteq S \cup \{v\}].$$

Algorithm 1 StochasticGreedyAJM (V, f, k)

```

1: initialize  $S = \emptyset, P = \emptyset$ 
2: for  $v \in V$  do
3:   with probability  $\frac{1}{2}, P \leftarrow P \cup v$ 
4: while  $|S| < k$  and  $\exists v \in P$  such that  $f(S \cup v) \geq f(S)$  do
5:    $v = \arg \max_{v \in P} f(S \cup v) - f(S)$ 
6:    $S \leftarrow S \cup v$ 
7: return  $S$ 

```

Now consider some $S' \supseteq S, v \notin S'$. The indicator function in the above sum counts edges where one of the two nodes is not contained in S , but both are contained in $S \cup \{v\}$. If S' extends S but does not contain v , then the summation for $h(v|S')$ can only include more nonzero terms than the summation for $h(v|S)$. Since each term is nonpositive, $h(v|S') \leq h(v|S)$. \square

In fact, h being monotone-decreasing and submodular is sufficient for the objective f to also be submodular:

PROPOSITION 4.2. *Whenever h is a monotone-decreasing submodular function, f is submodular.*

PROOF. Consider the marginal gain of adding a node v to a given seed set S :

$$f(v|S) = h(S \cup \{v\})a^T x_v + [h(S \cup \{v\}) - h(S)] \sum_{u \in S} a^T x_u$$

We prove that f is submodular by showing that each corresponding term in $f(v|S')$ can only decrease for all $S \subseteq S'$. The first term decreases since h is a submodular function as shown in the lemma. The second term, corresponding to the individual contribution of v , also decreases because h is monotonically decreasing. Thus $f(v|S') \leq f(v|S)$. \square

Having shown that f is submodular, a natural approach to seeding would use the greedy algorithm, giving a $1 - \frac{1}{e}$ approximation for the ICM and the LTM [8, 23, 25]. However, since f is non-monotone, this approach does not apply. Instead, we adopt the stochastic greedy algorithm proposed by Feldman, Harshaw, and Karbasi [11]. Algorithm 1 runs the normal greedy algorithm (lines 4-6) but only selects from a limited set of nodes P . Each node is included in P with probability $1/2$. This random removal reduces the chance that the greedy method will prematurely commit to a node that later become problematic due to non-monotonicity. Feldman et al. [11] show that this algorithm obtains a guaranteed $\frac{1}{4}$ -approximation to the optimal value and has excellent empirical performance. Our experiments follow their suggested strategy of running the algorithm several times (we both use 4).

4.1 Meta-Analysis of Influence Spread Metrics

The finding that the ICM is a poor predictor of node-level influence is dissonant with the fact that seeding algorithms based on the ICM have proved effective in the field [39]. After all, how can a seeding algorithm based on an inaccurate model of diffusion manage to nevertheless achieve a high level of influence spread? To address this seeming conflict, we confront a larger question about the prevailing methodology of the influence maximization problem. In this section,

we show that appealing solely to the magnitude of influence spread achieved is a fundamentally inconclusive method of determining whether a particular diffusion model underlies an observed instance of spread. This ambiguity is problematic when using the influence maximization framework to inform the seeding strategies of network interventions in sustainability domains. In many such cases, in addition to diffusing information generally, programs seek to target particular individuals or groups, and thus a model's ability to make node-level predictions is a valuable asset for real-world use cases.

Magnitude of Influence Spread Previous research comparing information diffusion model predictions to empirical results has tended to rely on metrics related to volume of spread—such as minimizing RMSE as a function of actual spread or recapitulating cascade sizes—to determine the fidelity of a model to ground truth processes [17]. However, we argue that one cannot extrapolate *processes* from such coarse-grain *outcomes*. The following experiments use three examples of physical, meso-scale networks: *Homeless*, a network of 142 nodes gathered via interviews with homeless youth, *India*, a household-level network gathered from a rural village in India [3], and *SBM*, a synthetic network of 200 nodes generated via the Stochastic Block Model, which replicates the community structure found in real social networks.

We evaluate how seed sets selected under one diffusion model perform in an influence maximization task under the other models. Figure 4a examines the consequences of model misspecification for influence maximization. We set the parameters equally across all networks—0.1 for propagation probabilities in the ICM and AJM and edge weights in the LTM. Each table entry shows the percentage of optimal influence spread obtained when a seed set selected according to the model on the column is assessed with the model on the row. For example, the cell (ICM, LTM) indicates that a seed set selected via the greedy algorithm for the LTM produced influence spread that was 99.8% optimal when diffusion actually occurred under the ICM. Given that all entries are greater than 90%, determining model fit by examining the magnitude of influence spread achieved under its seeding strategy leads to great ambiguity. Since all of these models result in high influence spread, *any* model could account for the “true” underlying process of diffusion.

One explanation for this phenomenon points to the community structure common in social networks. Algorithms for influence maximization under the ICM tend to distribute seed nodes across different communities to avoid the redundancy of seeding the same community multiple times. But seeding according to the AJM results in a similar recommendation to ensure diversity among seed nodes. Hence, attaining high influence spread is insufficient for identifying a model as the true diffusion mechanism. On the one hand, achieving comparable final magnitudes of influence spread is a handy property for influence maximization tasks, as it suggests that high-quality results are attainable even when the true model is uncertain. However, many important influence maximization tasks require a descriptively accurate diffusion model, not just one that works by coincidence.

Next, in Figure 4b, we show that common diffusion models are capable of reproducing *any* observed level of total influence spread. Each plot gives the fraction of the graph influenced by a random set of 10 seed nodes under the ICM and LTM as we vary a parameter for each model. For the ICM, we vary the propagation probability p .

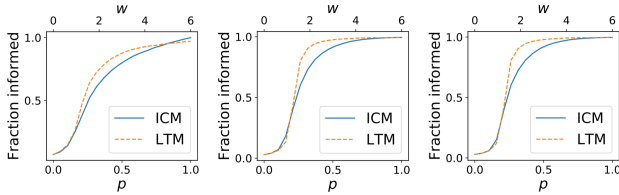
Figure 4

	ICM	LTM	AJM
ICM	100	99.8	98.6
LTM	99.6	100	98.8
AJM	97.4	96.1	100

	ICM	LTM	AJM
ICM	100	98.4	99.8
LTM	99.9	100	98.9
AJM	93.7	97.8	100

	ICM	LTM	AJM
ICM	100	98.7	99.4
LTM	99.3	100	99.8
AJM	96.3	93.9	100

(a) Percentage of optimal influence spread achieved when network is seeded according to the column model and evaluated according to the row model. Networks top to bottom: Homeless, India, SBM.



(b) Fraction of all nodes informed under ICM and LTM with varying parameters. Networks left to right: Homeless, India, SBM.

For the LTM, we assign each edge (u, v) a weight $\frac{w}{deg(v)}$ and vary w . Each line is an average of 30 draws of the random seeds. Under both models, any level of influence can be explained by a parameter choice in either of the models. We conclude that even if a given model exactly replicates the observed amount of influence spread in a network, this fact alone provides *no evidence* that the model truly describes the underlying diffusion process. Hence, we must use a finer-grained assessment such as node-level activations to produce accurate diffusion models.

Toward Node-Level Influence Spread The prevailing methodology’s blind spot to node-level information spread also engenders severe limitations in the actual deployment of the influence maximization problem in the real world. Seeding networks to maximize the total volume of influence spread is unproblematic when one is agnostic about *who* is influenced. But in many application domains, the objective for influence maximization is not to reach the greatest number possible out of *all* nodes, but rather to maximize the number reached in a given *subset*. In the Introduction, we referenced the variety of network interventions that seek to target particular individuals or sub-populations. Therefore, the total influence spread in the entire network is uninformative about the true objective. Achieving success in such a task requires a model which makes accurate node-level predictions, because this is the only way to find a seed set which influences the intended targets (while one can incorporate targeting into the ICM/LTM objective, this will only yield meaningful results if the model can accurately predict node-level activations). We have shown that the AJM significantly outperforms leading models in node-level predictions of who will be influenced by an intervention, making it much more suitable for applications.

5 DISCUSSION

A breadth of research has investigated information diffusion on on-line networks, but the problem of influence maximization remains under-explored on natural networks. By analyzing node-level data from a large-scale study of influence maximization on physical social networks, we show that neither of the prevailing models of information diffusion—the Independent Cascade and Linear Threshold—could account for the empirical findings. Even after fitting the best of these models to the data, they *perform worse than a random classifier* in predicting a node’s ultimate influence status. There is a long line of work in the AI literature on influence maximization under the ICM and LTM. Our results are significant because they open up a new setting for algorithm development and network modeling, moving beyond the prevailing models at least where physical networks are concerned.

We approached the shortcomings of the dominant models with an open mind to related research that may inform our understanding of information diffusion in this domain. Our proposed Activation Jump Model (AJM) draws from a lineage of work in multi-agent systems and social network theory that suggests that 1) social exchange need not only occur along network ties and 2) an individual’s influencing behavior is affected by her surrounding community. Of particular note, we model seed set structural diversity as conferring benefits to each individual seed node’s influencing level.

The AJM is a more inclusive model of diffusion and superior to leading models in its predictive prowess. When validated on three real-world networks with information diffusion data, the AJM issues predictions of node-level influence spread that *improve upon the best ICM and LTM predictions by 60% to 110%*. More generally, as the Activation Jump model relies on fewer network assumptions than standard graph approaches to information diffusion, we suggest that the framework may be of special interest to social network intervention programs where nodes have a high-level of familiarity with the rest of the network, and the act of social “jumping” is commonplace. Since the AJM is submodular and non-monotone, we adopt a seeding algorithm that achieves a $\frac{1}{4}$ -approximation to the optimal influence spread. Thus high-efficacy influence maximization under the AJM is computationally ready to be deployed in the real world.

It has long been accepted in the social sciences that the link between individual and group social behaviors is bidirectional [30]. A multi-agent team perspective is thus particularly suited to describe peer-to-peer information diffusion in the natural world, where a group’s social dynamic impacts how individual members will behave in spreading information. By explicitly modeling team-formation, a central component of many network interventions in the field, the AJM significantly updates the influence maximization problem for natural world settings. Its framework, with flexible activation level and landing distribution functional forms, also allows for contextually relevant information such as node-specific attributes like gender and ethnicity to be incorporated when deployed in real-world network applications.

Acknowledgments: This research was supported by MURI Grant W911NF-11-1- 0332, the California HIV/AIDS Research Program, and NSF Graduate Research Fellowships to Hu and Wilder.

REFERENCES

- [1] Noa Agmon and Peter Stone. 2012. Leading ad hoc agents in joint action settings with multiple teammates. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems, 341–348.
- [2] Tucker Balch. 2000. Hierarchic social entropy: An information theoretic measure of robot group diversity. *Autonomous robots* 8, 3 (2000), 209–238.
- [3] Abhijit Banerjee, Arun G Chandrasekhar, Esther Duflo, and Matthew O Jackson. 2013. The diffusion of microfinance. *Science* 341, 6144 (2013), 1236498.
- [4] Devon D Brewer. 2000. Forgetting in the recall-based elicitation of personal and social networks. *Social networks* 22, 1 (2000), 29–43.
- [5] Robert S Broadhead, Douglas D Heckathorn, David L Weakliem, Denise L Anthony, Heather Madray, Robert J Mills, and James Hughes. 1998. Harnessing peer networks as an instrument for AIDS prevention: results from a peer-driven intervention. *Public health reports* 113, Suppl 1 (1998), 42.
- [6] Ronald S Burt. 2009. *Structural holes: The social structure of competition*. Harvard University Press.
- [7] Carter T Butts. 2003. Network inference, error, and informant (in) accuracy: a Bayesian approach. *social networks* 25, 2 (2003), 103–140.
- [8] Wei Chen, Yajun Wang, and Siyu Yang. 2009. Efficient influence maximization in social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 199–208.
- [9] Morris H DeGroot. 1974. Reaching a consensus. *J. Amer. Statist. Assoc.* 69, 345 (1974), 118–121.
- [10] Tom Fawcett. 2006. An introduction to ROC analysis. *Pattern recognition letters* 27, 8 (2006), 861–874.
- [11] Moran Feldman, Christopher Harshaw, and Amin Karbasi. 2017. Greed is Good: Near-Optimal Submodular Maximization via Greedy Optimization. *arXiv preprint arXiv:1704.01652* (2017).
- [12] Samuel R Friedman, Melissa Bolyard, Pedro Mateu-Gelabert, Paula Goltzman, Maria Pia Pawlowicz, Dhan Zunino Singh, Graciela Touze, Diana Rossi, Carey Maslow, Milagros Sandoval, et al. 2007. Some data-driven reflections on priorities in AIDS network research. *AIDS and Behavior* 11, 5 (2007), 641.
- [13] M Gaston and M DesJardins. 2003. Team formation in complex networks. In *Proceedings of the 1st NAACSOS Conference*.
- [14] M Gaston, John Simmons, and M DesJardins. 2004. Adapting network structure for efficient team formation. In *Proceedings of the AAAI 2004 fall symposium on artificial multi-agent learning*.
- [15] Amer G Ghanem, Srinivasa Vedanarayanan, and Ali A Minai. 2012. Agents of influence in social networks. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems, 551–558.
- [16] Amit Goyal, Francesco Bonchi, and Laks VS Lakshmanan. 2010. Learning influence probabilities in social networks. In *Proceedings of the third ACM international conference on Web search and data mining*. ACM, 241–250.
- [17] Amit Goyal, Francesco Bonchi, and Laks VS Lakshmanan. 2011. A data-based approach to social influence maximization. *Proceedings of the VLDB Endowment* 5, 1 (2011), 73–84.
- [18] Mark Granovetter. 1978. Threshold models of collective behavior. *American journal of sociology* 83, 6 (1978), 1420–1443.
- [19] Mark S Granovetter. 1973. The strength of weak ties. *American journal of sociology* 78, 6 (1973), 1360–1380.
- [20] Barbara J Grosz, Sarit Kraus, Shavit Talman, Boaz Stossel, and Moti Havlin. 2004. The influence of social dependencies on decision-making: Initial investigations with a new game. In *Proceedings of the Third International Joint Conference on Autonomous Agents and Multiagent Systems-Volume 2*. IEEE Computer Society, 782–789.
- [21] Lu Hong and Scott E Page. 2004. Groups of diverse problem solvers can outperform groups of high-ability problem solvers. *Proceedings of the National Academy of Sciences of the United States of America* 101, 46 (2004), 16385–16389.
- [22] John Kalafat and Maurice Elias. 1994. An evaluation of a school-based suicide awareness intervention. *Suicide and Life-Threatening Behavior* 24, 3 (1994), 224–233.
- [23] David Kempe, Jon Kleinberg, and Éva Tardos. 2003. Maximizing the spread of influence through a social network. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 137–146.
- [24] David Krackhardt and Martin Kilduff. 1999. Whether close or far: Social distance effects on perceived balance in friendship networks. *Journal of personality and social psychology* 76, 5 (1999), 770.
- [25] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance. 2007. Cost-effective outbreak detection in networks. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 420–429.
- [26] Somchaya Liemhetcharut and Manuela Veloso. 2012. Modeling and learning synergy for team formation with heterogeneous agents. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 1*. International Foundation for Autonomous Agents and Multiagent Systems, 365–374.
- [27] Mahsa Maghami and Gita Sukthankar. 2012. Identifying influential agents for advertising in multi-agent markets. In *Proceedings of the 11th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*. International Foundation for Autonomous Agents and Multiagent Systems, 687–694.
- [28] Leandro Soriano Marcolino, Albert Xin Jiang, and Milind Tambe. 2013. Multi-agent team formation: diversity beats strength?. In *IJCAI*, Vol. 13. 279–285.
- [29] Peter V Marsden. 2005. Recent developments in network measurement. *Models and methods in social network analysis* 8 (2005), 30.
- [30] George Herbert Mead. 1934. *Mind, self and society*. Chicago University of Chicago Press.
- [31] Seth A Myers, Chenguang Zhu, and Jure Leskovec. 2012. Information diffusion and external influence in networks. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 33–41.
- [32] Foster J Provost, Tom Fawcett, Ron Kohavi, et al. 1998. The case against accuracy estimation for comparing induction algorithms.. In *ICML*, Vol. 98. 445–453.
- [33] Eric Rice, Anamika Barman-Adhikari, Norweeta G Milburn, and William Monro. 2012. Position-specific HIV risk in a large network of homeless youths. *American journal of public health* 102, 1 (2012), 141–147.
- [34] Thomas Schelling. 1978. *Micromotives and macrobehavior*. Nueva York. (1978).
- [35] Termeh Shafie. 2015. A multigraph approach to social network analysis. *Journal of Social Structure* 16 (2015), 0_1.
- [36] Johan Ugander, Lars Backstrom, Cameron Marlow, and Jon Kleinberg. 2012. Structural diversity in social contagion. *Proceedings of the National Academy of Sciences* 109, 16 (2012), 5962–5966.
- [37] Thomas W Valente. 2012. Network interventions. *Science* 337, 6090 (2012), 49–53.
- [38] Yongqing Wang, Hua-Wei Shen, Shenghua Liu, and Xue-Qi Cheng. 2013. Learning user-specific latent influence and susceptibility from information cascades. *arXiv preprint arXiv:1310.3911* (2013).
- [39] Amulya Yadav, Bryan Wilder, Eric Rice, Robin Petering, Jaih Craddock, Amanda Yoshioka-Maxwell, Mary Hemler, Laura Onasch-Vera, Milind Tambe, and Darlene Woo. 2017. Influence maximization in the field: The arduous journey from emerging to deployed application. In *Proceedings of the 16th Conference on Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents and Multiagent Systems, 150–158.
- [40] Ke Zhou, Hongyuan Zha, and Le Song. 2013. Learning triggering kernels for multi-dimensional Hawkes processes. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*. 1301–1309.