

THE DYNAMICS OF SOCIAL CONTAGION

A Dissertation

Presented to the Faculty of the Graduate School
of Cornell University

in Partial Fulfillment of the Requirements for the Degree of
Doctor of Philosophy

by

Vladimir Barash

August 2011

© 2011 Vladimir Barash
ALL RIGHTS RESERVED

THE DYNAMICS OF SOCIAL CONTAGION

Vladimir Barash, Ph.D.

Cornell University 2011

Social contagion is a subset of contagion which includes all social phenomena that can and do spread via social networks. The notion of how something becomes popular is very relevant to the concept of social contagion. Rumors, fads, and opinions can spread through social networks like wildfire, “infecting” individuals until they become the norm. This thesis investigates the dynamics of social contagion, employing a combination of formal analysis, simulation, and empirical data mining approaches to examine the processes whereby social contagion spreads throughout social networks. I introduce the concept of critical mass for a subclass of social contagion called complex contagion. This concept builds on earlier work to describe the nonlinear dynamics whereby most socially contagious phenomena infect very few people while a few become overwhelmingly popular. I also investigate socially contagious phenomena that arise when rational agents act under conditions of local information. Finally, I examine how my analytic work applies to a large dataset of empirical social contagion and draw implications for further research in the area.

BIOGRAPHICAL SKETCH

Vladimir Barash is a Ph.D. candidate in Information Science at Cornell University. At Cornell, Vladimir is a member of Prof. Michael Macy's lab, where he researches diffusion in social networks and social mechanisms as they apply to online communities. Prior to attending Cornell, Vladimir received his B.A. from Yale University. Vladimir was born in Russia, and currently lives in Ithaca, NY.

To my mother. To my friends. To all the people who have helped me be a better person.

ACKNOWLEDGEMENTS

Thanks goes first and foremost to my family, for supporting me for the last five years, in many many ways. Second, to my closest friends for being there for me through all my wanderings in and out of grad school. Thanks to Michael, the best advisor anyone could wish for, who was always ready to help when I needed help, but left me alone when I needed to be alone. Thanks to Jon and Claire, wonderful committee members who took time out of their very busy schedules to give me work and life advice on many occasions. Thanks to all my coworkers and colleagues at Cornell - especially to DanCo, who was always there when I needed to ruminate over beer. Last, but definitely not least, thanks to Scott Boorman at Yale and Eric Smith at Santa Fe and Duncan Watts at Yahoo! - the scientists who inspired me and pushed me as a wee young undergraduate to follow the path of Science.

TABLE OF CONTENTS

Biographical Sketch	iii
Dedication	iv
Acknowledgements	v
Table of Contents	vi
List of Tables	viii
List of Figures	ix
1 Introduction	1
1.1 Social contagion	1
1.2 The Structure of Social Networks	4
1.2.1 Network Properties	4
1.2.2 Network Models	7
1.3 Mechanisms of Social Contagion	10
1.3.1 Influence	11
1.3.2 Homophily	13
1.3.3 Environmental Effects	15
1.3.4 Social Identity	16
1.3.5 Local Information	19
1.3.6 Social Exclusion	23
1.4 Mathematical Models of Social Contagion	24
1.4.1 Node States	25
1.4.2 Success of a Contagious Phenomenon	25
1.4.3 Infection Functions	26
1.4.4 Models of Contagion	27
1.4.5 The SIR Model and Variants	27
1.4.6 The SIS Model	29
1.4.7 The SI Model	30
1.5 Empirical Analysis of Social contagion	31
2 Chapter 2: Local Information and Social Contagion	33
2.1 Background	34
2.2 General Model	34
2.3 Bounded Rationality and Myopia	35
2.4 Deviation from Optimal Behavior	37
2.5 Optimization of Local Information	38
2.6 Constraints on Optimization	39
2.7 Strategy Spaces and Agents	40
2.8 Network Structure	43
2.8.1 Poisson Random Graph	44
2.8.2 Small World Graph	45
2.8.3 Preferential Attachment Graph	45
2.8.4 Other Models	46

2.9	Analysis and Results	47
2.9.1	General Results	47
2.10	Poisson Random Graph	49
2.11	Small World graph	54
2.12	Preferential Attachment Graph	61
2.13	Simulations	67
2.14	Discussion	71
3	Chapter 3: Complex contagion and Critical Mass	74
3.1	Background	74
3.2	Model	77
3.3	Analysis	78
3.3.1	Small World Networks	80
3.3.2	Power Law Networks	83
3.3.3	Critical Behavior	86
3.3.4	Estimation of Function Behavior	89
3.3.5	Absolute and Relative Thresholds	94
3.4	Thresholds and Contagion Dynamics	96
3.5	Beyond the Threshold Model	100
3.6	Discussion and Conclusion	104
4	Chapter 4: Empirical Social contagion	108
4.1	Introduction	108
4.2	Theory	111
4.2.1	Threshold	111
4.2.2	Critical Mass	113
4.3	Data and Methods	116
4.3.1	Tag Coding	119
4.4	Results	119
4.4.1	Threshold	119
4.4.2	Critical Mass	122
4.5	Case Studies	126
4.6	Discussion and Conclusion	131
5	Chapter 5	135
5.1	Summary of Results	135
5.2	Broader Implications	138
5.3	Future Work	141
	Bibliography	143

LIST OF TABLES

LIST OF FIGURES

2.1	ρ_{crit} vs. size of infected cluster for Poisson Random Graph	68
2.2	ρ_{crit} vs. average bias for Poisson Random Graph	68
2.3	r vs. average bias for Unrewired Lattice Graph	69
2.4	r vs. average bias for Unrewired Lattice Graph without “shortcuts”	70
2.5	time to infect all nodes vs. ρ_{crit} for Barabasi-Albert graph	71
3.1	Inflection and drop-off points in the probability of \mathbf{a} rewired ties to an infected node as the level of infection increases on a rewired lattice, $\mathbf{a} = 2$, $\mathbf{k} = 48$, $\mathbf{p} = 0.1$, $\mathbf{N} = 40000$	90
3.2	Critical mass increases with threshold \mathbf{a} and decreases with perturbation \mathbf{p} on a rewired lattice, $\mathbf{k} = 48$, $\mathbf{N} = 40000$. Colors indicate critical mass from red ($\mathbf{CM} = 1$) to white ($\mathbf{CM} = \mathbf{N}$)	92
3.3	Inflection and drop-off points in the probability of \mathbf{a} rewired ties to an infected node as the level of infection increases on a power law network, $\mathbf{a} = 2$, $\mathbf{r} = 1$, $\mathbf{N} = 40000$, $\alpha = 2$	93
3.4	Critical mass increases with threshold \mathbf{a} and decreases with skewedness \mathbf{r} on a power law network, $\mathbf{N} = 40000$, $\alpha = 2$. Colors indicate critical mass from red ($\mathbf{CM} = 1$) to white ($\mathbf{CM} = \mathbf{N}$)	94
3.5	Contagion growth rate (black) and probability that a random node will be uninfected and have a random ties to infected nodes (red), for a regular lattice with $\mathbf{k} = 8$, $\mathbf{a} = 2$, $\mathbf{N} = 40000$, $\mathbf{p} = .1$	100
3.6	\mathbf{NIc} (logged) as a function of the log of the area on a perturbed lattice, $\mathbf{k} = 48$, $\mathbf{p} = .1$, $\mathbf{N} = 40000$, and $\mathbf{a} = 2$ (green) and $\mathbf{a} = 3$ (brown)	102
3.7	Expected number of adopters as a function of \mathbf{I} using only $\mathbf{f1(1)}$ (black), all of $\mathbf{f1}$ (red), and $\mathbf{f2}$ (green), $\mathbf{N} = 40000$, $\mathbf{k} = 48$, $\mathbf{p} = .1$	103
4.1	Confidence score distribution by threshold value for Flickr tags	120
4.2	Distribution of threshold value by manually labeled category for Flickr tags	121
4.3	Histogram of persistence parameters for Flickr tags	122
4.4	Contagion criticality $\rho(c)$ vs. number of adopters for Flickr tags	123
4.5	Contagion criticality $\rho(c)$ vs. threshold $a(c)$ for Flickr tags	124
4.6	Average tie range for ties through which contagion spreads vs. percentile of adopters	125
4.7	Threshold and criticality statistics for six contagion	127
4.8	Number of adopters and contagion growth rate by time for the iFlickr tag	128
4.9	Number of adopters and number of redundant ties from adopters to non-adopters by time for iFlickr tag	129
4.10	Average number of adopter friends by time for the iFlickr tag (brown - all tag adopters, blue - only tag adopters with one or more adopter friends)	130

4.11 Average number of adopter friends by time for six different tags 132

CHAPTER 1

INTRODUCTION

1.1 Social contagion

Social contagion is a subset of contagion which includes all social phenomena that can and do spread via social networks. The notion of how something becomes popular is very relevant to the concept of social contagion. Rumors, fads, and opinions can spread through social networks like wildfire, infecting individuals until they become the norm. What was originally a minority belief can become a dominant one as more and more individuals are exposed to said belief from their friends and choose to adopt it.

Products can also become socially contagious, via word-of-mouth marketing. In some cases, the success of a product can be entirely or mostly attributed to traditional marketing methods, such as TV or newspaper advertisements, which spread through informational networks (from TV stations or printing presses directly to the homes of consumers). In other cases, however, word-of-mouth recommendations by friends can account for many purchasing decisions.

Social movements are a third kind of social contagion. What is initially an unpopular or even marginalized cause can gather strong popular support as individuals convince their family, friends and acquaintances to join.

For consistency of terminology, I refer to social contagion as the set of all social phenomena that can and do spread via social networks. I refer to specific instances of social contagion (e.g. a specific product that spreads via viral

marketing) as socially contagious phenomena.

The mechanisms whereby social contagion spreads are the subject of active research in the computer, social and physical sciences. Understanding these mechanisms is key to the central research questions surrounding social contagion:

1. How far will a particular contagious phenomenon spread, i.e. how many individuals will it infect before it stops infecting new ones?
2. Given a particular networked population, what is the ideal seed for a socially contagious phenomenon, i.e. given the ability to infect some small subset of a networked population with this phenomenon, the choice of which particular subset will result in the greatest number of subsequent infected?
3. What are the dynamics of social contagion spread, i.e. given a particular networked population, will social contagion spread through certain parts of this population faster than they do through others?

This thesis will explore the three research questions using a combination of mathematical analysis, simulation, and empirical investigations of social contagion. The focus will be on the third research question, as the dynamics of contagion on a particular networked population inform both the ideal seed choice and the ultimate number of infected individuals. While this question has been received a lot of attention in previous work [72, 68], I explore the difference between the spread of social contagion and other contagion like information and disease through social networks. In particular, I introduce a concept of critical mass that is specific to a class of social contagion called complex contagion

[21]. Previous work [47] has examined the concept of critical mass in general, but the critical mass I propose is unique to complex contagion. Critical mass for complex contagion is tied to a phase transition in the range of ties through which complex contagion spreads, a phase transition that does not appear in the spread of non-social contagion like information and disease. My results show that the dynamics of social contagion are substantively different from the dynamics of non-social contagion, both in the case of simulated contagion and in the case of empirical contagion spreading on real-world networks.

The rest of this thesis is organized as follows. I begin with a background chapter on the structural properties of social networks and on the mechanisms whereby social contagion spreads through these networks. I also present a brief overview of the mathematical models of social contagion and recent empirical studies of social contagion, which will help frame the analysis in subsequent chapters. In Chapter 2, I examine the mechanism of local information and its effect on contagion dynamics, and prove several interesting theorems about agent decisions to adopt or not a contagious phenomenon spreading through different kinds of networks under conditions of local information. In Chapter 3, I focus on a particular class of social contagion known as complex contagion [21], and show that there exists a critical mass of adopters beyond which complex contagion is highly likely to spread throughout the entire network. Critical mass is particularly relevant as some real-world behaviors and products, such as social movements and virally marketed goods, may be examples of complex contagion. In Chapter 4, I examine a set of empirical contagious phenomena - tags in the Flickr photo sharing network - and attempt to apply the theoretical results from Chapter 3 to predicting their dynamics. I discover confirmation for the theory of complex contagion and provide evidence that some empirical con-

tagious phenomena reach critical mass. Furthermore, I observe that it may be possible to anticipate the critical mass point for these phenomena by analyzing the average estimated threshold of the contagious phenomenon from adoption statistics. I conclude with a discussion and high-level remarks in Chapter 5.

1.2 The Structure of Social Networks

I begin with a brief theoretical discussion of the structure of social networks. The simplest way to formally represent a network is a graph $G = (V, E)$, where V is the set of nodes or vertices, and E is the set of links or edges between the vertices. Members of V are usually represented with a unique identifier (name, address, etc.) and members of E are tuples of members of V . Some networks allow multiple edges, in which case the tuples do not have to be unique. Other networks allow self-edges, which means tuples of the form (i, i) are possible. These tuples can be ordered or unordered, depending on whether the edges represent asymmetric (citation of one paper by another) or symmetric (collaboration on a paper) relationships. If a graph has ordered tuples, it is called a directed graph, or digraph for short; otherwise, it is called an undirected graph.

1.2.1 Network Properties

Given this formalism, one of the simplest properties to describe is the degree of a node, and its macro-level equivalent, the degree distribution. In undirected graphs, the degree of a node i is simply the number of tuples (u, v) where $u = i$ or $v = i$. In a directed graphs, we consider the out-degree of i , which is the

number of (u, v) where $u = i$, and the in-degree of i , which is the number of (u, v) where $v = i$. The degree distribution of a graph is the function $P(x)$ which for any integer x gives the proportion of nodes i that have degree (in-degree, out-degree) equal to x . Naturally, $P(x)$ can take on any form, but by far the most common degree distribution in empirical networks is a power law, where $P(x) \approx x^{-\alpha}$, and α varies between 2 and 3 [3, 52]. This applies to both undirected and directed networks, where both in-degree and out-degree follow a power law [3, 52].

If we consider any two nodes i, j in a graph, we define a path between them as a sequence of connected nodes starting with i and ending with j . The geodesic distance or path length between i and j is the number of edges connecting the nodes in the shortest path between i and j (if no path between i and j exists, the geodesic distance is infinite). At the network level, researchers measure the mean geodesic, which is the arithmetic or harmonic mean path length of the entire network. Almost all empirical social networks have a mean geodesic on the order of $f(\log(N))$ where N is the number of nodes [72] and f is a polynomial function. Networks with path length polylogarithmic in N have been called “small worlds,” after the title of the Stanley Milgram experiment [48] that first discovered the phenomenon, and in popular literature this property is called “six degrees of separation,” as the median number of connections in Milgram’s original experiment was six.

We now consider, for some undirected graph G , any node i and the set of its neighbors J , the nodes that i is connected to by an edge. Given any two neighbors $j, k \in J$, we say the tuple (i, j, k) is a connected triple. Further, if $(j, k) \in E$, we say that (i, j, k) forms a triangle or transitive triple. Intuitively, the more tri-

angles there are in the graph, the more it resembles a giant super-dense cluster. This intuition leads to the network-level property of clustering coefficient CC , the fraction of triangles to all connected triples in G . First appearing in math and physics literature [13], this property has been studied in the sociology literature under the name fraction of transitive triples. Researchers have found many empirical networks that have high clustering coefficients independent of network size [3, 52].

A more local clustering coefficient has been proposed by Watts and Strogatz [72]. For any node i in undirected graph G , consider the number of ties t_{ij} present between j neighbors of i . The maximum number of such ties is $k(k-1)/2$, where k is the degree of i . Then the fraction $\frac{2|t_{ij}|}{k(k-1)}$ is defined as the Watts-Strogatz clustering coefficient of i , $CC_{WS}(i)$. The network-scale equivalent of this quantity, CC_{WS} is simply the average of $CC_{WS}(i)$ over all nodes. The Watts-Strogatz clustering coefficient has been studied in the sociology literature under the name “network density” [63] and found to differ from the fraction of transitive triples, as it weights the contributions of low-degree vertices more than the latter quantity. Still, empirical networks have high, size-independent values of network density, as with their fractions of transitive triples.

At the highest level, a network consists of one or more components. We say nodes i and j are in the same component (strongly connected component for directed networks) if j can be reached from i following the edges of the network. Otherwise, i and j are in different components. The number and size distributions of network components can give insight into the cohesiveness of its structure, by which I mean the fraction of all nodes that are in the giant component. Studies of empirical networks [3] have almost uniformly found the presence

of a giant component, with 80% or more of the nodes as members, and an exponential distribution of smaller components, indicating that many networked systems are extremely cohesive. Related studies in network resilience [4, 16] have looked at the distribution of component number and sizes in a network as nodes are gradually removed from it. Empirical networks turn out to be robust against the removal of random nodes, the giant component persisting with as many as 70% of the nodes removed. At the same time, empirical networks are vulnerable to the removal of high-degree nodes, and fall apart into a set of tiny components with $O(\log(n))$ members after just 2-3% of the highest-degree nodes have been removed.

1.2.2 Network Models

Early work in network analysis focused on recording properties of real networks and constructing network models, often described as simple algorithms that would give rise to networks with the same properties. Though initial work on network models was done in the 1950s, the field has exploded with publications starting around a decade ago. There are now several comprehensive reviews of network models [3, 52]. We begin the section with the simplest network models, called random graph models, and then discuss generalized random graphs, small world networks, and network growth models.

Poisson Random Graphs

The first research into network modeling was done independently by Rapoport [58] and Erdos and Renyi [31] who proposed the very simple Random Graph

model, which continues to be used as a comparative tool for other, more complex models to this day. The Random Graph model consists of placing N nodes and then adding some $n < N(N-1)/2$ edges between them. Alternatively, we can say that every edge between the N nodes is present with a probability $p = n/N$.

Random graphs present an extremely simple model for network structure with an interesting transition point from a disconnected assembly of small, tree-like components to a giant component with many loops and a small diameter. Properties of the random graph model above the transition point are evocative of most empirical networks, which also feature a giant component with a small diameter. While the model fails to capture many properties of empirical networks, such as power-law degree distribution and high clustering coefficient, it is a good baseline for studying network phenomena and evaluating other, more complex models. In particular, the random graph model is an excellent example of how global network structure arises even in the complete absence of local network structure, a phenomenon that suggests that networks are complex artifacts with multiple levels of organization.

Small World Graphs

The Small World graph model interpolates between two undirected graphs: an Erdos-Renyi random graph and a lattice (the simplest lattice is a ring of nodes, where every node has exactly two neighbors). The Erdos-Renyi random graph has a small mean geodesic but low clustering coefficient, the ring lattice has a large mean geodesic but high clustering coefficient, so interpolations of the two graphs may exist that have both small mean geodesic and high clustering coefficient. Indeed, by starting with a ring lattice and randomly rewiring

edges, Watts and Strogatz [72] were able to achieve an intermediate regime where mean geodesic falls off drastically while clustering coefficient remains high. Some researchers refer to this regime as the “small world region,” [72] between the regularized world of the lattice and the random environment that emerges when too many edges are rewired.

Researchers often make use of the small world model when studying social networks, as it captures two critical properties of these systems. The robust local connections of the lattice resemble small groups that abound in real social networks (neighborhoods, groups of friends, reading circles), while the shortcuts resemble the “weak ties” [34] observed empirically by sociologists since the 1970s (acquaintances, distant relatives). There is a lot of evidence that weak ties are a critical property of social networks, allowing information to spread quickly across many nodes and giving individuals access to resources (e.g. potential employers) that they would be hard-pressed to discover in their local social circles. Weak ties are also an important catalyst of social integration, bringing together groups of different people that would otherwise never interact socially. Naturally, the small world model has become a useful tool for studies that investigate weak ties, the spread of information and search in social networks.

Preferential Attachment Graphs

The Preferential Attachment graph model was formulated independently, first by Price [26], then by Barabasi and Albert [12]. The critical principle of the model states that when v joins a network and an edge is created between v and w , the likelihood of some specific node w' being picked as w is proportional to its degree $k(w')$. This phenomenon was first studied by Herbert Simon [64] in the

1950s. Price applied Simon's methods to social networks, and Barabasi and Albert coined the term "preferential attachment" (PA). The latter authors in [12, 4] made two critical findings about PA: 1) the model generates networks with a power law degree distribution, and 2) without either the network growth mechanism or the preferential attachment principle, the degree distribution does not follow a power law. These findings suggest that preferential attachment captures a key element of the growth dynamics of empirical networks.

In addition to having a power-law degree distribution, networks generated by the PA model have path lengths logarithmic in the number of nodes, and commensurate with path lengths found in empirical networks of similar size (as opposed to those generated by generalized random graph models). The clustering coefficient of PA networks is dependent on network size, following $C \sim N^{-.75}$. In contrast, the clustering coefficient of empirical networks seems to be independent of network size.

1.3 Mechanisms of Social Contagion

We now turn to a background review of the mechanisms of social contagion. These mechanisms are the various social forces that make products, rumors and social movements spread from person to person. Influence is perhaps the mechanism most often associated with social contagion, but there are others. In this section, I will describe in detail influence, as well as local information, social identity, and social exclusion mechanisms that cause individuals to adopt things their friends have already adopted. In addition, I describe homophily and environmental factors as selection mechanisms that create effects that may on the

surface resemble the diffusion of social contagion but have distinct underlying causes.

1.3.1 Influence

Perhaps the most well-studied mechanism of social contagion is influence. Alters, individually or as a group, influence the ego to become more like them. As a result, over time the ego's actions or beliefs grow to reflect the actions and beliefs of her alters. Key to the notion of influence is that of culture, a set of "beliefs, attitudes and behaviors" [9] espoused by an individual. A common representation of this set is a collection of features (e.g. Language, Religion, Style of Dress) with a set of traits corresponding to each feature (Language can be French, English, and so on). The cultural profile of an individual is the set of traits, one per feature, that she most closely identifies with. The central hypothesis of influence-based theories is, then, that the ego's cultural profile changes over time to more closely resemble the profiles of her alters. In the context of social contagion, influence implies that if one or more of the alters espouses a belief or adopts a product, the ego will be more likely to do so.

An alternative form of influence is negative influence, where alters influence the ego to become less like them. This phenomenon has received much less coverage in the literature, but it does appear in cases of teenage rebellion and, in general, contentious social relationships. Negative influence does not play a strong role in social contagion, as individuals would be influenced to not adopt the contagious phenomenon their friends adopt.

It is important to note that even though influence is a prominent mecha-

nism of social contagion, it is often triggered by other mechanisms like social identity, so it is important to study those mechanisms when examining the role of influence in the spread of contagion. I describe these mechanisms in subsequent sections. Similarly, the actual manifestations of influence as a mechanism vary in social networks, from persuasion to enforcement of conformity through threat of sanction to memetic spread (e.g. role modeling).

Influence accounts for the formation of many social structures. Local communities arise in the course of personal identification with the in-group along key cultural dimensions like religion and political beliefs (e.g. [44, 18], which may arise in the course of local interactions, e.g. [24, 57]. At larger scales, there is some evidence that nations are more likely to form when their people have shared meanings and interlocking habits of communication [27, 28]. Similar patterns occur in transnational integration and succession conflicts. Other studies have linked influence to the spread of social norms [43, 67, 8], the spread of knowledge [19], the diffusion of innovations [59, 51] and the establishment of technical standards [61, 10].

There are two main kinds of influence - social and interpersonal. Interpersonal influence accounts for the ability of high-status, impressive individuals to dictate the behavior of their friends. Sometimes, all it takes is one popular alter to purchase a product, to instantly compel ego to do the same. Social influence accounts for the ability of groups to put social pressure upon an individual. The impact of social influence on behavior changes is cumulative with the number of alters who join in the behavior.

Influence is often triggered by other mechanisms of social contagion as, for instance, by social identity. It is important to distinguish between general the-

ories of social and interpersonal influence as a mechanism of contagion (which focus on the process of influence rather than its underlying causes) and theories of social contagion that include social influence as part of a more complex mechanism (and focus on the underlying causes that trigger social influence and thus enable the spread of social contagion). In this thesis, we focus on the latter set of theories, as the process of influence has been well investigated in the literature [9]. Furthermore, a number of empirical studies [25] have looked at the process of influence as an empirical phenomenon but few have delved further into its underlying causes.

1.3.2 Homophily

The concept of selection is, in the most general sense, that an individual's attributes may select for or against ties between that individual and others. A more specific version of selection is the notion of homophily, which states that an individual is more likely to develop social relationships with people who are like him or her in some respect. Homophily can be formalized in the same cultural profile framework as selection: over time, the set of social network alters for a particular individual changes, based on that individual's cultural profile. New ties appear to alters with similar cultural profiles, old ties to alters with different cultural profiles disappear. At any given point in time, homophily predicts that ties between people with similar cultural profiles are much more likely to exist than ties between people with different cultural profiles. In this section, we focus on homophily in particular, as research of this phenomenon has formed the bulk of selection-oriented studies.

The classic social science study on homophily is Lazarsfeld and Merton's [40] observation of friendships in Hilltown and Craftown. The researchers drew on earlier theoretical work, including anthropological studies of homogamy (homophily in marriage), to investigate the frequency of ties between similar people in empirical data. They also applied the now-popular quotation "birds of a feather flock together" to describe this phenomenon. Even before Lazarsfeld and Merton, small-scale studies [15, 45] showed substantial homophily by demographic characteristics like age, sex and education. Later work has looked at homophily at larger scales, in schools [29], communities [69] or even the US population as a whole [17].

Homophily is not a mechanism of social contagion, as it does not act as a force that helps a behavior spread from one individual to another in a social network. However, as a recent study [7] points out, homophily produces patterns that closely resemble the spread of social contagion. If two individuals *A* and *B* share a tie because they are similar, they may decide to adopt a behavior independently of one another, but at different times. The pattern of the *A* adopting, then *B*, will resemble that of the contagious phenomenon spreading from the *A* to *B* (possibly as a result of influence of *A* upon *B*). For example, two people *A* and *B* who are both avid runners may become friends through their common love of running. When choosing a running shoe, *A* and *B* may choose the same brand (even if neither is exposed to advertisements from this brand - see below), because they have similar tastes in shoes. Due to chance, *A* buys the shoe first, and then the pattern of *A* buying the shoe, then her friend *B* buying the shoe suggests that *A* influenced *B* when in fact it was homophily that induced both to buy the same brand.

1.3.3 Environmental Effects

Environmental effects are another mechanism that creates effects that resemble diffusion on the surface, when in fact no diffusion takes place. For instance, consider a rain cloud that moves over a park. The movement of the rain cloud is followed by a sequence of opening umbrellas in the crowd of park visitors, but umbrella opening is not a socially contagious phenomenon, and no diffusion takes place. Environmental effects are especially important when analyzing large-scale adoption data: consider the spread of topical messages (tweets, instant messages) through a social media network (Twitter, the Google Chat network). In the case of topics like “breakfast”, it is easy to mistake an environmental effect (rapid adoption starting with accounts in the East Coast of the United States and moving to accounts in the West Coast of the United States) for a diffusion pattern.

Environmental effects can interact with diffusion and homophily. Consider the above example of two friends A and B who have bonded over their love of running and have similar tastes in shoes. Then both A and B might be similarly affected by the same advertisement for a particular shoe brand and both decide to buy the shoe. Again, let’s assume that A buys the shoe first. Then the temporal pattern of A buying the shoe then B buying the shoe might be mistaken for A ’s influence on B , and the shoe brand for a complex contagion, when in fact the mechanism of brand adoption was an interaction of homophily (mutual love of running, similar tastes in shoes) and environmental effects (exposure to the same advertisement).

Such patterns are extremely difficult to distinguish from true patterns of contagion diffusion. The key difference between the two processes is that, for a

contagious phenomenon to be said to be adopted by *B* after it is adopted by *A*, *B* needs to be exposed to *A*'s having adopted the phenomenon (the exposure need not be deliberate on *A*'s part). If *B* sees *A* behaving in a particular way and then decides to behave in the same way (for whatever reason, so long as it is related to witnessing *A*'s behavior), then the phenomenon spreads from *A* to *B*. If, however, *A* and *B* adopt the behavior independently, without seeing each other engage in it, but rather due to a predisposition they share to engage in said behavior (a sharing that implies some similarity between *A* and *B*), or due to a predisposition they share to react similarly to a third party *C*, then we cannot say a contagious phenomenon has spread from one to the other. Empirical data analysis studies, even those with extremely rich data, often do not have access to recorded exposure events. Experimental studies can control for exposure, but lack the scale that data analysis can achieve. In this thesis, we mention existing approaches for distinguishing true diffusion (whether as a result of influence or other causes) from independent adoption events due to selection, whether the selection is for similar environmental effects or similar friends (selection of alters). We also indicate opportunities for future research in distinguishing diffusion mechanisms from selection mechanisms in large-scale settings.

1.3.4 Social Identity

When analyzing social contagion, it is important to examine the perspective of adopter identities, both at the individual and at the group level. Social Categorization Theory [66] posits that identity is a continuum of self-categorizations. At one extreme, there is an emphasis on interpersonal differences and identity manifests as a set of personal traits that define an individual as set apart from

all other individuals. At the other extreme, there is an emphasis on intragroup similarities and intergroup differences, so identity manifests as a set of social categories that define an individual as a prototypical member of some group in contrast to all members of different groups. A particular individual's identity constantly shifts along this continuum, as she perceives social categories as more or less salient to her beliefs and actions. These shifts also correspond in changing the salience of interpersonal vs. social influence (e.g. in dyadic vs. triangular relationships) and as such are highly relevant to the effectiveness of influence as an adoption mechanism.

Turner and Oakes show that, as social categories become more salient, the individual assumes a social identity that marks her as a member of some group and places her at odds with members of different groups (here we assume for simplicity that an individual can only belong to one group at a time). In these conditions, social influence is highly relevant. The social identity creates a perception of similarity between the individual and in-group members, which leads the individual to expect a level of agreement between herself and the in-group on all stimuli. A socially contagious phenomenon adopted by the rest of the group but not by the individual in question is a new stimulus which violates this assumption, creating a cognitive dissonance: in-group members are similar to the individual, similar people in the same setting should behave similarly [66], but there is a difference of behavior between the in-group and the individual with respect to contagion adoption. In response to this cognitive dissonance, the individual can:

1. "[Ignore] the cognitive dissonance and maintain the status quo
2. "[Attribute] the disagreement to perceived relevant differences between

self and others,”

3. “[Attribute] the disagreement to perceived relevant differences in the stimulus situation,”
4. “[Engage in] mutual social influence [with the in-group members] to produce agreement.”

Option 4 is often equivalent to the individual adopting the contagious phenomenon as a result of social influence, especially if most of the in-group has adopted the contagious phenomenon. The particular option chosen depends on the individual, group and behavior in question. In this thesis, we are obviously interested in option 4, but make sure to control for possible alternatives such as options 1, 2 and 3 when investigating the effect of social identity on contagion adoption in an empirical setting. Thus social identity acts as a mechanism for the spread of social contagion.

One specific aspect of social identity relevant to the spread of social contagion is the effect of triadic closure between in-group members on contagion adoption. If two adopter friends of a non-adopter share a tie (forming a closed triad between themselves and the adopter), all three are more likely to belong to the same group and share a social identity. As Turner and Oakes have argued, this makes the adopters more likely to exert social influence upon the non-adopter by virtue of their shared social identity, to the extent that a 2-1 imbalance in adoption status within the closed triad group makes adoption the dominant behavior in this instance of social contagion. So social identity supports the hypothesis that non-adopters will be more likely to adopt a socially contagious phenomenon if their adopter friends are friends with each other. Empirical studies that demonstrate the positive effect of triadic closure on adop-

tion probability [11] are consistent with social identity. In this thesis, we investigate the role social identity plays in this specific hypothesis via empirical data analysis and experimental testing.

1.3.5 Local Information

Local information is a mechanism of social contagion, that, unlike influence or social identity, does not affect the likelihood of infection, but, instead, affects the likelihood of exposure. The origins of local information are in bounded rationality theories [64] that attempt to explain why agents deviate from optimal strategies given by rational choice theory.

One assumption of rational choice theory is that people have all the information they need to make any sort of decision. In the real world, people often lack critical pieces of information when making decisions, and so their decisions might deviate from those predicted by rational choice theory.

A specific kind of incomplete information we are interested in is local information about the structure of connections between agents in a network-embedded population. This concept appears in the economics literature on network games [32] where “a player’s well-being depends on own actions as well as actions taken by his or her neighbors.” Network games in general are complex to analyze; the authors note that even very simple games played on networks have multiple equilibria with a wide range of possible outcomes. Introducing incomplete information in global games tends to reduce equilibrium multiplicity [50], but the precise kind of incomplete information to introduce is somewhat arbitrary. Galeotti et al. argue that network games suggest two

natural sets of incomplete information: information about the identity of future neighbors and about the number of neighbors they will have. We focus on the second of these suggestions and leverage it to define local information as follows: the local information of an agent playing a network game is the states and connections between her neighbors.

Local information relates to the spread of social contagion via the concepts of network externality [38]. In the economics literature a network externality is the relationship between the number n of people who buy a product and the utility of buying the same product for the $n + 1$ st person¹ In the context of contagion, agents deciding whether to adopt a particular contagious phenomenon may base their decision on a utility function that is strictly increasing in the number of adopters. Note that the concept of network externality would apply even for contagious phenomena that have no monetary value, such as joining a social movement, so long as the decision to adopt a contagious phenomenon can be expressed as a utility function that compares cost and benefit, and either benefit goes up or cost goes down in the number of current adopters.

In the case of perfect information, the agent will decide whether to adopt the contagious phenomenon based on how many total adopters there are up to that point. In reality, agents are often blind to such global behavioral patterns on the population scale, but pay a lot of attention to similar patterns at the local neighborhood scale: for instance, they may not know how many people in the world own an iPod, but they know very well how many of their friends do. These local patterns, in turn, may represent a biased sample of the entire population, and thus distort the agent's perception of global trends. To the extent this distortion

¹We focus here on positive network externalities, where the relationship is monotonically increasing.

occurs, the agent may be said to be deviating from perfectly rational behavior due to the effect of local information.

Early studies of contagion adoption [34] assume that everyone in the population of potential adopters knows everyone else, so local information is equivalent to perfect information. As a result, the spread of a contagious phenomenon in Granovetter's model is irreversible given a super-critical distribution of thresholds. In contrast, Morris [49] and Centola and Macy [21] explore the dynamics of threshold-based social contagion that can spread only through local interactions, and conclude such contagion spreads most effectively when network structure at the local level is robust (many redundant ties). Centola and Macy's complex contagion may begin to spread throughout a network, but stop after a few iterations due to local information effects embedded in the network structure.

Models of local information do not have to rely entirely on local network neighborhoods. In network-embedded models, agents may iteratively query ever more distant nodes to get better information before adopting a particular product. Tie strength is another relevant dimension, as highly local interactions may arise when agents rely only on their closest friends (an even smaller set than all their network neighbors!) when making an adoption decision.

In Chapter 2, we explore in detail a model where agents have to rely on local information to make their decision about adopting a particular contagious phenomenon, but may query more distant nodes, which may improve the information they receive, but incur a (constant) cost. We constrain the agents to make their decisions prior to the spread of the contagious phenomenon. This constraint mimics several real-life situations such as groups fighting "misinfor-

mation contagion.” Consider the case of a community leader who wants to prevent drugs from entering her community but knows that individuals will make decisions based on the actions of their friends rather than some rational cost-benefit analysis. Should the community leader enforce “quarantine”² and attempt to limit external influences on the community, cutting off outside ties as much as possible; or should she encourage more outreach, hoping that as the network size of community members increases, individuals will have ever more heterogeneous networks, decreasing the likelihood that a new behavior (such as drug use) will be dominant among their neighbors? Similar decisions may arise for strategists trying to prevent contagion like climate change denial (or climate change belief), and so on, from spreading to a target community. Note that in these cases, it may not be possible (or prudent) to prevent the target community from ever adopting a particular contagion, but it may be possible to delay (or accelerate) adoption to avoid bias from local information. In other words, consider some global utility u_g of adopting a contagious phenomenon c , expressed as a fraction of adopters in the entire population. We would like to investigate optimal behavior with respect to adopting c exactly when u_g of all agents have adopted, not sooner or later.

Our motivation in Chapter 2 is based on the intuition of Galeotti et al. [32] that: “when players have limited information about the network they are unable to condition their behavior on its fine details and this leads to a significant simplification and sharpening of equilibrium predictions.” Indeed, we find that it is not necessary to calculate explicit equilibrium conditions but rather to focus on the general dynamics of a contagious phenomenon spreading through a network. These dynamics naturally suggest strategies for a set of nodes wishing to

²The notion of quarantine is captured by the SIR model of contagion, see below.

avoid adopting the contagious phenomenon too early (or too late) due to local information effects.

The model in Chapter 2 resembles the “local knowledge” model of Michael Chwe [22], who investigates games with participation thresholds based on the decisions of network neighbors to participate or not. Chwe also looks at the possibility that nodes query a mixture of local neighbors and distant nodes by introducing a parameter into his model that governs the frequency of random connections in the network. I discuss the differences between Chwe’s and my model in more detail in Chapter 2.

1.3.6 Social Exclusion

Two previously described mechanisms of social contagion (influence and social identity) enable the diffusion thereof by encouraging non-adopters to adopt. Social exclusion acts in a diametrically opposite way: it enables diffusion of social contagion by discouraging non-adopters from continuing to not adopt.

Social exclusion is a concept from social politics that encapsulates the denial of access to economic, political or cultural systems which determine the social integration of a person in society. [70]. It is important to note that this term is often used in context of social change and in general perceived as a negative aspect of social order. In this work, we focus on the narrow, technical problem of how social exclusion interacts with diffusion of social contagion. Therefore, we examine processes of social exclusion in the abstract, outside of a social change framework. Nevertheless, it is important to note that the results of our work may have important implications for empirical cases of social exclusion, and

may cast light on certain aspects of social politics.

From a technical perspective, social exclusion has two important properties: dynamics and locality. The first property is noted by [70] who describe social exclusion as “the dynamic process of being shut out, fully or partially, from any of the social, economic, political or cultural systems which determine the social integration of a person in society.” The second property is referenced in [46] who define social exclusion as a multi-dimensional process that find[s] a spatial manifestation in particular neighborhoods. Social exclusion is a process that happens over time rather than a static property of groups, and is situated in local neighborhoods as opposed to global networks.

1.4 Mathematical Models of Social Contagion

Contagion is any behavior, message, product or organism that spreads throughout a population. Contagious phenomena were first studied in the context of epidemiology, as organisms (like bacteria) that spread from one member of a population to another. Much more recently, interdisciplinary work in physics, computer science and sociology has studied contagion more broadly as behaviors (e.g. smoking), messages (rumors), and products (gadgets) that spread throughout populations. The body of work on contagion is voluminous and diverse, and the associated terminology is not uniform, but a few key terms are used commonly. I outline these terms first to inform the rest of the discussion.

1.4.1 Node States

Contagious phenomena are most often characterized by a change in the state of individuals within a larger population, from uninfected to infected. Other states, like resistant or vaccinated, are possible but the binary distinction between uninfected and infected individuals is key to representing contagion. In all contagion models, uninfected individuals can become infected. In some contagion models, infected individuals can become uninfected again (people can get over a disease, give up a product, or change their behavior).

1.4.2 Success of a Contagious Phenomenon

Researchers often speak of the success or failure of a contagious phenomenon. Again, several definitions of contagion success exist, but the common one is that the success of a contagious phenomenon is the proportion of infected nodes at equilibrium, i.e. when nodes no longer change state from infected to uninfected (or vice versa). In models where the infection frequency at equilibrium is binary (either close to 100% or close to 0% nodes infected), the term takeover is used to describe situations where the contagious phenomenon infects nearly all of the nodes at equilibrium, and the term failure is used to describe situations where the phenomenon takes over almost none of the nodes at equilibrium. In the rest of this work, we will use the terms takeover and failure for binary infection frequencies at equilibrium and the term success for continuous infection frequencies at equilibrium.

1.4.3 Infection Functions

The third key element of any contagion model is the process whereby uninfected nodes become infected. This process is often represented by a function from a set of individual, local, and/or population-level attributes to a probability of infection. There are two major types of infection functions: deterministic or threshold functions, where the probability of infection is always either 1 (when input attributes take on certain ranges of values) or 0 (otherwise); and probabilistic functions, where the output probability of infection is a continuous variable. Mixtures of deterministic and probabilistic functions are possible, but rarely used.

Infection functions also vary by their input variables: in fully-mixed contagion models, any individual in a population can be infected by any other individual at any time. Therefore, infection functions take as input only individual (e.g. node state) or population-level (e.g. frequency of infected individuals) variables. In partially-mixed contagion models, individuals can only be infected by a subset of the population, identified by geographic proximity, social ties, etc. In partially-mixed contagion models, infection functions take as input individual and local-level (e.g. frequency of infected individuals among friends) variables, but most often ignore population-level variables.

Some contagion models have recovery functions. These represent the opposite process of an infection function - the process whereby infected nodes become uninfected. Recovery functions can also be threshold or probabilistic, but we do not study them in great detail here.

1.4.4 Models of Contagion

Finally, there are a number of contagion models in the literature. All of these models attempt to describe the dynamics of contagion as it spreads from some initial subset of the population (often called the seed) to other individuals, until equilibrium in the state of individuals (infected vs. non-infected) is reached. We present a brief overview of contagion models in the scientific literature, starting with early epidemiological models and moving on to later computer science, physics and social science applications.

1.4.5 The SIR Model and Variants

The earliest well-known model of contagion is the so-called SIR model [1, 5, 6, 36]. This model was first formulated (though never published) by Lowell Reed and Wade Hampton Frost in the 1920s. It divides a population into 3 classes: susceptible to some contagious phenomenon (S), infective (I) and Recovered (R). The original work simply assumes that any susceptible individual has a uniform probability β per unit time of catching the disease from an infective one, and that infective individuals recover and become immune at some constant rate γ . Recovered individuals can never become infected again. The distribution of the three classes over time is then governed by a system of differential equations. The SIR model is a fully-mixed model with a probabilistic infection function and a probabilistic recovery function.

The SIR model exhibits critical behavior in its parameters. If β is low and γ high, the phenomenon almost never spreads throughout the population. Conversely, if β is high and γ low, the phenomenon almost always takes over the

entire population. In between these extremes, the model exhibits a non-linear transition in the values of β and γ : if β is below a certain critical value, most nodes remain uninfected, but if β exceeds that critical value, the phenomenon rapidly takes over the entire population. Similarly, a decrease in γ below some critical value leads to takeover whereas populations with higher values of γ remain contagion-free. This property of the SIR model is reproduced in most models of contagion we study below and is the focus of much interest in research on the diffusion of contagion.

Grassberger [35] extended this fully mixed model onto a n -dimensional lattice network (making it a partially-mixed model, again with a probabilistic infection and recovery functions), where a susceptible node s can become infective with probability β if and only if it has an infective neighbor, and any infective node becomes recovered after a fixed length of time γ passes since its infection. The same paper shows that the networked version of the SIR model can be mapped exactly onto results from a field of study called percolation theory. Percolation theory arose in physics as the study of critical phenomena in current flow. The major result of the field as it applies to network structure is that, starting from a graph of many small connected components (sets of nodes that are connected to each other but to no other node in the graph), and adding edges at random to this graph, one arrives at a critical phase where the addition of just a few edges instantly connects the many small components into one giant component. Grassberger's work showed that if the graph formed by connections between all the infected nodes percolates (has a single giant connected component), the partially-mixed SIR model predicts contagion takeover of the entire graph, while if the same graph does not percolate (has many small connected components), the partially-mixed SIR model predicts contagion failure.

There are a few other variants of the SIR model that are worth mentioning. Pastor-Satorras and Vespignani [54, 55] analyze the SIR model on networks with power law degree distributions and show that for all non-zero β , the contagious phenomenon takes over the network as long as the power law exponent is less than 3. An interesting variant of the SIR model includes “vaccination,” the removal from a network of some particular set of vertices prior to the contagious phenomenon simulation. Pastor-Satorras and Vespignani [56] discuss this variant at length, noting that removal of the highest-degree vertices makes the network disconnected and prevents contagion takeover (a result that is inspired by studies of network resilience). Cohen et al. [23] investigate the vaccination problem in a context where information about the network is limited and propose an interesting workaround, to follow random edges in the network as those are most likely to lead to nodes with the highest degree.

1.4.6 The SIS Model

An alternative to the SIR model is the SIS model, where the dynamics are as in SIR, but infective nodes turn back to susceptible ones and never fully recover. This model also exhibits critical behavior in the values of β and γ , transitioning from contagion failure to contagion takeover. SIS cannot be solved analytically as SIR but Pastor-Satorras and Vespignani [54, 55] give a detailed investigation of the model on a class of simulated networks known as configuration model networks[53], and show that the contagious phenomenon persists in such networks (some non-zero fraction of the population is infected) for all non-zero β values in this network class. A recent study by Leskovec et al. [42] uses an SIS model to replicate data on information diffusion through a network of blogs.

1.4.7 The SI Model

A third diffusion model is SI, which dates back to Mark Granovetter's work on threshold models of collective behavior [34]. Granovetter considered a fully mixed population of susceptible individuals with different thresholds. An individual i becomes infective if t_i or more other infective individuals are currently in the population, where t_i is i 's infection threshold, and all infective individuals remain so forever. Granovetter showed that this model also exhibits critical behavior: below a critical density of infective nodes, the model always resulted in contagion failure, but at or above this density, the model always resulted in contagion takeover. Granovetter's model is fully-mixed with a deterministic infection function and no recovery function.

Morris [49] constructed a partially-mixed (network embedded) model based on Granovetter's work with the restriction that a node becomes infective if some fraction k of its neighbors are infective, so the infection threshold is uniform for all nodes. Morris investigated the maximum k for which diffusion of a contagious phenomenon can occur on an arbitrary size network, found an upper bound $k \leq 1/2$ and observed that actual maximum k are close to $1/2$ when the network resembles a lattice.

Watts [71] studied the same model in configuration model networks under the name "cascading failure". Watts was able to identify a high-connectivity regime in which cascades were very rare but occasionally took over the entire network (as in Granovetter's work), and suggested that in the latter regime, the ideal cascade seed neighbored a lot of average-degree nodes, as their thresholds could be satisfied with fewer neighbors than the thresholds of very high-degree individuals.

Centola and Macy [21] focus on the Morris version of the SI model and study the dynamics of this model on unrewired lattices and Small World networks. Centola and Macy examine both the relative thresholds studied by Granovetter and Morris, and absolute thresholds where a node is guaranteed to adopt if some raw number a of its neighbors have adopted. They find that contagion of any threshold (up to the maximum dictated by Morris' $k \leq 1/2$) spread in a similar way on an unrewired lattice network. However, in a Small World network, a difference between "simple" contagion of threshold 1 (or $1/z$ where z is the number of neighbors), and "complex" contagion of threshold > 1 (or $> 1/z$) emerges.³ Simple contagion spreads quickly through Small World networks, whereas complex contagion spreads more slowly or not at all, depending on the actual threshold value and the level of rewiring used to generate the network. We discuss Centola and Macy's work in much more detail in Chapter 3, as it forms the basis for our analysis of complex contagion and discovery that these phenomena have a critical mass of adopters.

1.5 Empirical Analysis of Social contagion

Empirical analysis of social contagion is a relatively recent field. Early investigations of diffusion have focused on online communities, which are seen as potential targets for viral marketing. A virally marketed product would spread on the network of connections in an online community much like a contagious phenomenon, so work on LiveJournal [11] and blogs [2, 39] has been met with a lot of interest. The first of these papers is especially interesting, as it finds

³The terms "simple" and "complex" contagion were coined by Centola and Macy in the same paper

that local network structure correlates heavily with probabilities of joining a group in the LiveJournal site. Specifically, individuals with many ties between their friends (high local clustering coefficient) are more likely to join groups than those with few ties between their friends. Other studies [41] look product recommendations in an online site, and make several empirical observations relevant to cascade models, including: the second recommendation of a product heavily increases the probability of buying, but further recommendations do not; the number of recommendations exchanged between two people does not positively influence the probability of buying; and there are categories of products for which recommendations are very effective, including professional/technical and personal/leisure items (but not works of fiction).

More recently, empirical analyses of contagion have been extended to voice mail products spreading through phone networks [30] and hashtags spreading through Twitter [60]. We anticipate many more studies in this area as large-scale data recording the adoption of behaviors and products online become available to the research community.

CHAPTER 2

CHAPTER 2: LOCAL INFORMATION AND SOCIAL CONTAGION

In this chapter, I explore in detail a model where agents have to rely on local information to make their decision about adopting a particular socially contagious phenomenon, but may query more distant nodes, which may improve the information they receive, but incur a (constant) cost. I constrain the agents to make their decisions prior to the spread of the contagious phenomenon. This mimics several real-life situations such as groups fighting “misinformation contagion.” Consider the case of a community leader who wants to prevent drugs from entering her community but knows that individuals will make decisions based on the actions of their friends rather than some rational cost-benefit analysis. Should the community leader enforce “quarantine” and attempt to limit external influences on the community, cutting off outside ties as much as possible? Or should she encourage more outreach, hoping that as the network size of community members increases, individuals will have ever more heterogeneous networks, decreasing the likelihood that a new behavior (such as drug use) will be dominant among their neighbors? Similar decisions may arise for strategists trying to prevent contagious phenomena like climate change denial (or climate change belief), and so on, from spreading to a target community. Note that in these cases, it may not be possible (or prudent) to prevent the target community from ever adopting a particular contagious phenomenon, but it may be possible to delay (or accelerate) adoption to avoid bias from local information. In other words, consider some global utility u_g of adopting a contagious phenomenon c , expressed as a fraction of adopters in the entire population. I would like to investigate optimal behavior with respect to adopting c exactly when u_g of all agents have adopted, not sooner or later.

2.1 Background

the contagious phenomenon model I use assume that agents are rational (that is, they will behave to optimize some utility function $u(d)$ where $d \in \{\text{adopt}, \text{not-adopt}\}$, but are limited to local information about the decisions of their neighbors, as opposed to global information about the decisions of everyone in the network. This model resembles the “local knowledge” model of Chwe [22]. Chwe also provides a detailed analysis that demonstrates the effect of local information on decision-making under conditions of different network structure and threshold. My work differs from Chwe’s in two important aspects. First, I do not model the contagious phenomenon problem as a simultaneous game, so agents in my model never make decisions based on what their neighbors might do, but only based on what their neighbors have done so far. Second, in my model individuals can change their neighborhoods so as to adopt as close as possible to the time step (if any) when they would have adopted given global information.

2.2 General Model

Consider a graph G where the vertices are abstract agents, and the edges are social ties between the agents. A contagious phenomenon C diffuses along this graph. Each agent $a \in A \equiv V(G)$ has an adoption state $S(a)$ where $S(a) = 1$ if the agent has adopted C and $S(a) = 0$ otherwise. Let us assume that the utility of adopting C is a step function in the density of adopters:

$$\begin{aligned}
u(\text{adopt}(C)) &= 1 \text{ if } D(C) > \rho_{crit} \text{ 0 o.w.} \\
u(\text{non-adopt}(C)) &= 1 - u(\text{adopt}(C))
\end{aligned}$$

where $D(C)$ is the density of adopters, the number of adopters $N(C)$ divided by the total size of the population P . Here ρ_{crit} is the critical density value above which adoption becomes the optimal strategy. At or below ρ_{crit} , non-adoption is the optimal strategy. Agents with perfect information will always adopt when $D(C) > \rho_{crit}$ and not adopt when $D(C) \leq \rho_{crit}$. From this model it follows that any contagious phenomenon C on G will not spread unless seeded with $\lceil \rho_{crit} P \rceil$ adopters, and will spread instantly if that condition is met.

2.3 Bounded Rationality and Myopia

We now introduce bounded rationality into the decision function. We will assume that agents are myopic, that is, they can see the state $T(a)$ only for $a \in L(a) \subset A$. Beyond $L(a)$, the agent is only aware of the existence of other agents, but not of their states nor of their network connections. We will investigate three scales of myopia:

- M1: $L(a)$ consists of one neighbor a' of a
- M2: $L(a)$ consists of all neighbors a' of a
- M3: $L(a)$ consists of all neighbors a' and all neighbors-of-neighbors a'' of a .

Under each level of myopia, the agent must make her decision solely based on the states of other agents in $L(a)$. To represent this limited decision-making, we introduce a new utility function u_L :

$$\begin{aligned} u_L(a, \text{adopt}(C)) &= 1 \text{ if } D_L(a, C) > \rho_{crit} \text{ 0 o.w.} \\ u_L(a, \text{non-adopt}(C)) &= 1 - u_L(a, \text{adopt}(C)) \end{aligned}$$

where L is the set of all local information sets $L(a)$ and $D_L(a, C)$ is the local adopter density of a , that is, the number of adopters in $L(a)$, $N_L(a, C)$, divided by $|L(a)|$. Myopic agents will adopt if $D_L(a, C) > \rho_{crit}$ and not adopt otherwise. For notational purposes, we consider u_L to be the adoption utility function with respect to L .

Under conditions of bounded rationality, far fewer than $\lceil \rho_{crit} P \rceil$ seeds are necessary for the contagious phenomenon to spread to the entire population. Consider the following scenario:

- A graph G represented by a sequence $\langle a_1, a_2, \dots, a_n \rangle$ such that each a_i is connected only to the preceding a_{i-1} (if exists) and the following a_{i+1} (if exists).
- M1 for this graph: a_{i-1} (if exists).
- A contagious phenomenon C with $\rho_{crit} = 1/2 - \epsilon$ for some small $\epsilon > 0$
- seeds for the contagious phenomenon: a_1 and a_2 .

Under these conditions, for each scale of myopia as defined above, a_3 will adopt C , then a_4 , and so on until the contagious phenomenon has spread to every

agent in the graph. Note that under conditions of bounded rationality, global adoption rate does not instantly change from 0 to 1 (as it did for perfect information), but goes up over time. We can define an artificial timeline for the adoption process by running a simulation where at t_0 the scenario is set up as above and a_1 and a_2 are seeded as adopters, and at t_1 and each timestep thereafter, each agent simultaneously makes a decision about whether or not to adopt C . Then at timestep t_1 node a_3 will adopt, at timestep t_2 node a_4 will adopt until at timestep t_{n-2} node a_n adopts and the entire set A has adopted the contagious phenomenon.

2.4 Deviation from Optimal Behavior

The scenario above shows how agents can act in a suboptimal manner due to the effects of local information. Let's assume $n = 20$, $P_{crit} \equiv \lceil \rho_{crit} P \rceil = 10$. Then, during a run of the simulation described above, agent a_3 acts in a suboptimal manner at times t_1 through t_9 : at these timesteps, she decides to adopt whatever the scale of myopia chosen, because her local decision function indicates that adoption is optimal ($D_L(a, C) > \rho_{crit}$), but the true decision function evaluated at these points indicates non-adoption as the optimal behavior ($D(C) \leq \rho_{crit}$). Similarly, agent a_{20} acts in a suboptimal manner at times t_{10} through t_{17} : at these timesteps, she decides to adopt because her local decision function indicates that non-adoption is optimal, but the true decision function evaluated at these points indicates adoption as the optimal behavior. In fact, only agent a_{12} acts optimally at all times because she decides to adopt starting precisely at the timestep when it becomes true that $D(C) > \rho_{crit}$.

2.5 Optimization of Local Information

The simulation example in the previous section suggests a general formulation of the local information problem: given a contagious phenomenon C spreading on a graph G of myopic agents, is it possible to predict which agents have such local information sets $L(a)$ that they will adopt C if and only if it is optimal to do so? A related problem is, given a subset of all agents on the above graph, is it possible to change their local information sets $L(a)$ so that they will adopt C if and only if it is optimal to do so? Again, optimality here is defined according to the true decision function, which mandates adoption if and only if the $u(\text{adopt}(C)) > u(\text{non-adopt}(C))$ and non-adoption otherwise.

These formulations, while ideal from a rational choice perspective, are too restrictive to be applicable to empirical contagion, where adoption events may be influenced by external factors and therefore perfect information about the utility of adoption cannot be captured in our model. However, we can examine the related question of whether certain local information sets can enable agents to act optimally as often as possible, controlling for external factors. This question can be restated as an optimization problem of local information:

Given a contagious phenomenon C spreading on a graph G of myopic agents and a subset $S(A)$ of all agents in the graph, is it possible to change the local information sets $L(a)$ for $a \in S(A)$ so that these agents act optimally as often as possible?

We can formalize this problem as the search for a transformation $OptLocal$ of the local information sets $L(a)$ for agents $a \in S(A)$ that generates new local informa-

tion sets for these agents that minimize the amount of time each agent spends acting suboptimally:

$$OptLocal(S(A)) = \underset{f(L)}{\operatorname{argmin}} \sum_{a \in S(A)} \sum_{t \in T} u(d_t(a)) - u_{f(L)}(a, d_t(a))$$

where $f(L)$ produces a new set of information sets L' that consists of the original information sets $L(a)$ for all agents $a \notin S(A)$ and new information sets $L'(a)$ for all agents $a \in S(A)$. The next step is to determine what sets L' produce the least bias and what transformations f generate these sets.

2.6 Constraints on Optimization

Without further constraints, the optimal transformation f is one that maximizes the size of local information sets of all agents in $S(A)$. To avoid this trivial solution, we impose a constant cost c for each bitwise difference between $L(a)$ and $L'(a)$ for each agent. To calculate this cost, first we order all the agents $a \in A$ to create a sequence $\langle a_1 \dots a_n \rangle$ where $n = |A|$. Now we can represent each local information set $L(a)$ as a binary vector $V(L(a))$:

$$V(L(a)) = \langle i : i = 1 \text{ if } a_i \in L(a) \text{ } 0 \text{ o.w. } \rangle$$

now the cost of transforming one local information set is simply:

$$Cost(L'(a)) = c * \|V(L(a)) - V(L'(a))\|$$

we can use this cost function to modify *OptLocal* to *OptLocalCost*, the cost-constrained optimal transformation of L :

$$OptLocalCost(S(A)) = \underset{f(L)}{\operatorname{argmin}} \sum_{a \in S(A)} Bias(f(L(a))) * Cost(f(L(a)))$$

where:

$$Bias(f(L(a))) = \sum_{t \in T} u(d_t(a)) - u_{f(L)}(a, d_t(a))$$

The optimization problem as I have set it up is linearly separable at the level of agents, that is, the choices of new local information set by one agent do not affect the choices of the other agents. This separability only holds under a further assumption, that the choice by agent a of agent b to be in $L(a)$ is a directed relationship (so a is not automatically in $L(b)$). I begin examining the problem under this assumption, and leave the implications of lifting this assumption for future work.

The new optimization problem to find $OptLocalCost(S(A))$ is non-trivial, and in fact the space of possible local information sets $L'(a)$ for a particular agent a is exponential: $|\{L'(a)\}| = 2^{|A|}$. I now narrow down his space by considering possible strategies for optimization.

2.7 Strategy Spaces and Agents

It is important to remember that, even though I have formulated local information as an optimization problem, standard quantitative techniques for solving

such problems may not be applicable to the model.

Many techniques for solving optimization problems rely on learning or iterative improvement. Such techniques, however, require global information: the agent must know how well she did at any given time in order to decide whether it was better or worse than previous times. For example, in the case of supervised learning, the agent can consult a teacher, who would tell the agent how often she behaved optimally vs. sub-optimally during the contagious phenomenon adoption process. The model I have set up purposefully avoids all sources of global information: a key hurdle for the agents to overcome is the lack of any reference point about the states of all but a few other agents. If a teacher were available, she could simply tell the agent everyone's state as they adopted, granting her perfect information which would automatically lead to optimal behavior.

Furthermore, these techniques require the contagious phenomenon process to happen multiple times, so that the agent has many opportunities to pick a local information set. In the model as set up, the contagious phenomenon process happens only once and is irreversible: both u and u_L are step functions, so all non-seed agents will not adopt until adopter density reaches a critical value, at which point they will adopt forever. Some empirical contagious processes (like getting a body piercing or joining a protest) behave according to this model, while others do not. This opens up an opportunity to introduce alternative models of local information, but for the purposes of my analysis I will focus on the model as presented.

Finally, it is possible to create highly tailored strategies that allow for optimal behavior. For instance, in the sequence graph described in Section 2.3, let's

consider as $S(A)$ the sequence of agents $\langle a_{13} \dots a_{20} \rangle$. These agents could change their local information sets to $\{a_{11}\}$, and then all would adopt precisely when it became optimal to do so. However, these tailored solutions again require global information: the agents have to be aware of the full set of connections in G , and of the seed nodes. Such level of knowledge is outside the scope of the model.

Absent iterative improvement and tailored strategies, what sorts of strategies are available to agents? To answer this, let's consider the knowledge agent a has about other agents. At a minimum, a knows well any of the agents in $L(a)$, and so can choose them to be in $L'(a)$. However, a may also know about the behavior of other agents not in $L(a)$: in empirical social networks, some very popular individuals, like celebrities, have their behavior broadcast via information streams, word-of-mouth, and other channels. As a result, many people may be aware of a celebrity's adoption state for particular contagious phenomenon, such as whether she endorses a specific product, without knowing her directly. I integrate the notion of popularity into the myopic model as follows: a can pick at random any agent b outside of $L(a)$ to be in $L'(a)$, with probability of picking some particular b proportional to b 's popularity, which I define as her degree.

The highly local nature of information available to agents in the myopic version of the model greatly reduces their strategy space. Only two strategies are available for agent a : pick an agent in $L(a)$, or pick a random agent not in $L(a)$ with probability proportional to that agent's degree. Using these two strategies, the agent makes up her new local information set $L'(a)$. I reflect this in the formal definition of the optimization problem by further constraining $f(L)$:

$$f(L) = L' = \{L(a) : a \notin S(A)\} \cup \{L'(a) : a \in S(A)\}$$

$$L'(a) = a_l^* a_g^* : a_l \in L(a), a_g \notin L(a) \wedge P(a_g) \propto k(a_g)$$

where $P(a_g)$ is the probability of picking a particular agent a_g not in $L(a)$ and the * symbol implies zero or more of the previous items, as in regular expressions. We can now write down the full set of equations that describes the cost- and agent-awareness- constrained optimization of the set of local information sets to optimize adoption behavior:

$$OptLocalCost(S(A)) = \underset{f(L)}{\operatorname{argmin}} \sum_{a \in S(A)} Bias(f(L(a))) * Cost(f(L(a))) \quad (2.1)$$

$$Bias(f(L(a))) = \sum_{t \in T} u(d_t(a)) - u_{f(L)}(a, d_t(a)) \quad (2.2)$$

$$f(L) = L' = \{L(a) : a \notin S(A)\} \cup \{L'(a) : a \in S(A)\} \quad (2.3)$$

$$L'(a) = a_l^* a_g^* : a_l \in L(a), a_g \notin L(a) \wedge P(a_g) \propto k(a_g) \quad (2.4)$$

2.8 Network Structure

So far, we have considered the problem of local information set optimization in the abstract. In particular instances of this optimization problem, the structure of connections between agents has a big effect on the flow of contagion through the network, and thus, on the optimality of different local information sets. As the example in section 2.3 shows, it is easy to predict the flow path of a

contagious phenomenon on particular network structures. However, the analysis of the flow of a particular contagious phenomenon on a particular network may not generalize to other networks, other contagious phenomena and other instances of the local information set optimization problem. Instead, I analyze contagion flow on several *network models*, or broad classes of network structures meant to replicate empirical networks (human society, the World Wide Web, social media). This section briefly describes these models.

2.8.1 Poisson Random Graph

The simplest network model is that of a Poisson Random Graph (cite - Erdos/Renyi), which is initialized as a set of N nodes with no connections between them. Subsequently, each pair of nodes is connected with probability p . The parameters p and N fully describe the model. A well-known property of Poisson Random Graphs is a phase transition around $pN = 1$. For $pN < 1$, the graph is a set of tiny components with no connections between them. For $pN > 1$, a giant component emerges. Poisson Random Graphs have low diameter like empirical networks, but are poor candidates of empirical social structures in other respects: they have a Poisson degree distribution (instead of a power law or a log normal), and a relatively low clustering coefficient. Nevertheless, the simplicity of construction for a Poisson Random Graph makes it a good starting point for studying the structure and dynamics of networks from an analytic and simulation perspective.

2.8.2 Small World Graph

The Small World Graph (cite - Watts/Strogatz) is initialized as an n -dimensional lattice with N nodes. Each node in the lattice has k neighbors. After initialization, ties are randomly added (cite - Watts) or rewired (cite - Maslov / Snep-pen) to create shortcuts in the lattice. The lattice is a degree-regular graph with high diameter and high clustering coefficient. The added or rewired ties creates shortcuts, which decreases diameter but also decreases clustering coefficient. Watts and Strogatz showed that, for a certain range in the density of shortcut ties, diameter decreases dramatically (to the levels of a random graph), but clustering coefficient remains high. Lattices with shortcut density in this range are known as small world graphs. Small world graphs are good models of many empirical social structures as they combine dense local structure (high clustering coefficient) with high connectivity (low diameter) that many researchers say is indicative of human social networks. For instance, many human social networks exhibit dense local structure around circles of friends, neighborhoods, or co-workers; at the same time, “weak” ties through distant acquaintances make these networks highly connected in what is known as the Six Degrees of Separation phenomenon in popular culture. The drawback of small world graphs is their uniform (for rewired shortcuts) or Poisson (for added shortcuts) degree distribution, which is not representative of many empirical networks.

2.8.3 Preferential Attachment Graph

Preferential attachment graphs (cite - Barabasi/Albert) are initialized with one or a small number of nodes with no connections between them. Subsequently,

the model adds new nodes to the graph one at a time. When a new node i appears, it immediately makes a number a of new connections to existing nodes. For each node j that is already not a neighbor of i , the probability of an $i - j$ tie is proportional to j 's degree. By this mechanism, nodes j that already have high degree are ever more likely to have even higher degree as new nodes i appear. This "rich-get-richer" effect on degree accumulation is known as preferential attachment in network analysis, and results in a power-law degree distribution that is representative of many empirical networks. Preferential attachment graphs also have low diameter, but have low clustering coefficient as well, lacking the dense local structure of small world graphs.

2.8.4 Other Models

Other models like Forest Fire(cite - Kleinberg) attempt to capture different aspects of the structure and dynamics of empirical networks. Some models explicitly seek to combine high clustering coefficient with power law degree distribution and low diameter. Other models focus on the dynamics of tie formation and attempt to replicate empirical tie formation processes. For the purposes of the local information model, I will focus on the three models presented above, as they cover a broad spectrum of network structures. My analysis of random, small world, and preferential attachment graphs should serve as a foundation on which future work can examine contagion flow and local information in more complex network models.

2.9 Analysis and Results

I now turn to analysis of the optimal local information set problem on various network structures. Before discussing particular network structures, it is important to consider the problem in a general way, to help frame the structure-specific analyses. Broadly, we are interested in making nodes in $S(A)$ becoming infected by some contagious phenomenon C that starts from some seed set $E(A)$ and spreads throughout the network, under very specific circumstances: that is, as close to as possible to the moment when C has infected exactly P_{crit} nodes.

2.9.1 General Results

Theorem 2.9.1. *Consider an instance of the local information set optimization problem as defined above, where $|E(A)|$ is $O(P_{crit})$. Then the optimal sampling strategy for nodes in $S(A)$ is to sample entirely from the set $L(a)$, assuming the graph is random, or heavily from the set $L'(a) = a_g^*$ otherwise. In particular, for degree-regular networks the optimal sampling strategy is to sample entirely from the set $L'(a) = a_g^*$, whereas for degree-skewed networks the optimal sampling strategy is to sample from the set $L'(a) = a_g[\rho_{crit}]a_l[1 - \rho_{crit}]$, in other words, to include nodes a_g with proportion ρ_{crit} in $L'(a)$ and include nodes a_l with proportion $1 - \rho_{crit}$ in $L'(a)$.*

Proof. In the case where $|E(A)|$ is $O(P_{crit})$, the optimal point for adopting the contagious phenomenon may be at the start of its diffusion, or very shortly after, so the goal for all nodes in $S(A)$ is to become infected right away. So the optimal strategy for all nodes in $S(A)$ is to construct a local information set that consists of at least ρ_{crit} infected nodes.

In the case of a Poisson Random Graph, the ties of any node represent a random sample of the network, so ρ_{crit} of them are expected to be in $E(A)$. In this case, the optimal strategy is to stick with the local information set $L(a)$ corresponding to the node's network neighborhood.

In the case where the network is degree-regular (such as a small-world graph), then sampling nodes from a_g^* yields a true random sample of the network, which, as for the Poisson Random Graph, is likely to contain ρ_{crit} nodes that are in $E(A)$.

Conversely, assuming the network is not heavily rewired, sampling nodes from a_l^* is a wasteful strategy: for nodes that are in $E(A)$ already, this strategy will with high likelihood produce a sample where nodes in $E(A)$ are over-represented, whereas for nodes that are not in $E(A)$ this strategy will produce a sample where nodes in $E(A)$ are under-represented. Since nodes have no way of knowing in advance whether they will be in $E(A)$ or not, the optimal strategy remains to sample entirely from a_g^* . It is possible to rely more heavily on a_l^* for networks that are heavily rewired and thus resembling random graphs, but it is also impossible to tell the rewiring level of the network without global information about it, and even when the network is completely rewired, sampling from a_g^* will continue to result in an optimal local information sets given the constraints of the optimization problem.

Thirdly, in the case where the network has a skewed degree distribution (such as a preferential attachment network), sampling from a_g^* will lead to a sample of high-degree nodes. These nodes are not likely to be in $E(A)$ (since there are very few of them in the network), but they are highly likely to be infected first once the diffusion process starts (since most ties from all over the

network, $E(A)$ inclusive, lead to them). Accordingly, making sure that the local information set has nodes from a_g with frequency ρ_{crit} will guarantee that a adopts at the first time step. \square

Theorem 2.9.2. *Consider an instance of the local information set optimization problem as defined above, where $|E(A)| \ll P_{crit}$ and the contagious phenomenon never reaches P_{crit} nodes. Then the optimal sampling strategy for nodes in $S(A)$ is to sample from $L(a)$.*

Proof. In the case when $E(A)$ is much smaller than P_{crit} . In this case, if the contagious phenomenon never reaches P_{crit} nodes then the goal is for no node in $S(A)$ to ever become infected. In particular, if the contagious phenomenon does not spread very far beyond $E(A)$, then the vast majority of the nodes in the network, and so, in $S(A)$, will behave optimally using only their local information sets L . \square

The third, intermediate case, is when $E(A)$ is much smaller than P_{crit} but the contagious phenomenon may eventually reach P_{crit} nodes. This is the most interesting case and the one we will explore in the following sections.

2.10 Poisson Random Graph

Consider a Poisson Random Graph with N nodes and probability of connection between two nodes p and some small set of seeds $E(A)$ with sum degree m , with the assumption that $|E(A)| \ll P_{crit}$. Let us further assume that p is $O(\frac{1}{N})$ so a giant component has formed, but not much larger, so the graph does not form a clique [52]. Then I claim that:

Theorem 2.10.1. For $\rho_{crit} \leq \frac{1}{N_p}$, contagion spreads throughout the entire network in a logarithmic number of time steps, so the average number of sub-optimal time steps for nodes in $S(A)$ is minimal given original local information set $L(a)$.

For $\frac{1}{N_p} < \rho_{crit}$, contagion dies out in a logarithmic number of time steps, so only $O(|E(A)|/N)|S(A)|$ nodes in $S(A)$ should behave suboptimally.

Proof. The key to this proof lies in the highly random structure of the PRG, which leaves very little room for redundant ties that would allow high-threshold contagion to spread.

In the case where $\rho_{crit} \leq \frac{1}{N_p}$, contagion can spread through just one tie between an infected and an uninfected node. This means that starting from $E(A)$, contagion will infect all the neighbors of $E(A)$, then all their neighbors, and so on until it takes over the entire network. In the extreme case where there are no redundant ties, contagion takes over the entire network in $O(\log_m(N))$ steps, so on average no node in $S(A)$ will behave suboptimally for more than a logarithmic number of time steps in the size of the network, even if the strategy chosen is to rely entirely on the original local information set $L(a)$. The number of redundant ties at each step “slows down” contagion, and so increases the time that nodes in the network spend behaving suboptimally if they rely only on $L(a)$.

In the case where $\frac{1}{N_p} < \rho_{crit}$, contagion must spread through at least two ties between infected nodes and an uninfected node. This means that starting from $E(A)$, contagion will infect only the neighbors of nodes in $E(A)$ who have a minimum of two ties to nodes in that set. If this set is smaller than $E(A)$, then the next set will be even smaller (due to the homogeneous structure of the random

network), and so on, so instead of spreading to the entire network in a logarithmic time, contagion will “die out” in $\log_{m_{red}}(1/m)$ time steps, where m_{red} is the fraction of redundant to all ties. In the case where there are no redundant ties, contagion will only infect $|E(A)|$ nodes, so only a very small fraction of nodes in $S(A)$ is ever likely to become infected, with the rest behaving optimally (never becoming infected). More redundant ties will increase m_{red} , slow down the die-out time, and increases the number of nodes in $S(A)$ that behave suboptimally.

We now prove a lemma that provides an asymptotic bound for the fraction of m ties coming out of $E(A)$ that are redundant, for a Poisson Random Graph. This bound is very close to 0 for small m (which we assume at the beginning of this theorem), so we can assume the cases with no redundant ties described above hold in close approximation for this family of graphs, and the theorem is proved □

Lemma 2.10.2. *Consider a Poisson Random Graph defined as above in Theorem 8.3. Then for some set $E(A)$ of nodes with sum degree m , the total number of redundant ties among these m is $O((m/N)^2)$ which is close to 0 given the assumptions laid out in Theorem 8.3.*

Proof. Each of the m distinct ties coming out of $E(A)$ targets a node at random with uniform probability. The resulting number of n distinct nodes targeted by these m ties is, therefore, equal to the number of n distinct elements that results from a sampling m times from a population of N with uniform probability. This quantity is given by Tillé(2006) as:

$$N - \frac{(N-1)^m N!}{N^m (N-1)!}$$

which reduces to:

$$N \left(1 - \left[\frac{N-1}{N} \right]^m \right) \quad (2.5)$$

Focusing on the inner term, we have:

$$\left[\frac{N-1}{N} \right]^m = \left[1 - \frac{1}{N} \right]^m = \left[1 + \frac{-1}{N} \right]^m =$$

by binomial expansion:

$$= \sum_{k=0}^m \binom{m}{k} \left[\frac{-1}{N} \right]^k = \binom{m}{0} 1 + \binom{m}{1} \frac{-1}{N} + \binom{m}{2} \frac{1}{N^2} + \dots + \binom{m}{m} \left[\frac{-1}{N} \right]^m$$

This series S has the property that, for any $k = 0, k \leq m$, the $k + 1$ st element is smaller in magnitude and opposite in sign to the k th element. The sign opposition comes from the -1 in the power term of the series. The magnitude difference comes from the fact that the $k + 1$ st element is $O([m/N]^k)$, which decreases in k since $m < N$.

This property implies that the first few terms will dominate the series. In particular, we can establish bounds of the series with the first two partial sums: 1 and $1 - m/N$. Every subsequent term will alternatively drive the series closer to 1 and to $1 - m/N$, by an ever-decreasing degree, so the final sum will always stay within those bounds. Also note that, by the same property, the final sum will be much closer to $1 - m/N$ than to 1 . Accordingly, we can approximate the inner term as follows:

$$\begin{aligned} \left[\frac{N-1}{N} \right]^m &\geq 1 - \frac{m}{N} \\ &\approx 1 - \frac{m - \epsilon}{N} \end{aligned}$$

we can now rewrite Equation 5 as:

$$N \left(1 - \left[\frac{N-1}{N} \right]^m \right) \approx N \left(1 - \left[1 - \frac{m - \epsilon}{N} \right] \right) \quad (2.6)$$

$$\approx m - \epsilon \quad (2.7)$$

What does Equation 7 tell us? Instead of targeting m distinct nodes, m random ties in a PRG target some slightly smaller number $m - \epsilon$ nodes. In other words, $m - \epsilon$ ties target distinct nodes in the network, and the remaining ϵ ties are redundant.

We can approximate the magnitude of ϵ by taking the difference between the second and the third partial sums of S which is equal to:

$$\epsilon \approx \frac{m(m-1)}{2N^2} \text{ which is } O\left(\left[\frac{m}{N}\right]^2\right)$$

Since we assume p is $O(\frac{1}{N})$ and $|E(A)| \ll P_{crit}$, we have:

$$m = |E(A)|pN \approx |E(A)| \ll P_{crit} < N$$

and so $\epsilon \approx 0$ for large N .

□

2.11 Small World graph

Consider a Small World graph with N nodes where each node has k edges and p is the probability of any edge being rewired. Such graphs can often be represented as rewired dimensional lattices. In principle, the lattice can be of a dimension d , but this is not a parameter frequently used in Small World graph analysis, so for simplicity we keep $d = 1$. I first examine the simplest case where $p = 0$, where the network is an unrewired lattice. I present a lemma that gives the baseline bias for this case, i.e. the bias when the local information set of every a in $S(A)$ consists entirely of $L(a)$, a 's network neighbors.

Lemma 2.11.1. *For an unrewired lattice graph G with each node having k edges, and a contagious phenomenon seeded with one node e and its network neighborhood, so $E(A) = e \cup Nbrs(e)$, if the contagious phenomenon reaches at least P_{crit} nodes, the bias associated with local information set $L(a)$ for node a is:*

$$Bias(L(a)) = \rho_{crit} |D(a, E(A))k - P_{crit}|$$

where $D(a, E(A))$ is the graph distance between a and the seed cluster $E(A)$.

Proof. The case where contagion does not reach P_{crit} nodes is covered in subsection 8.1, so we are interested in the alternative - the case where contagion does reach P_{crit} nodes. Recall that in this case the bias for node a is given as the difference in time steps between the point when contagion infects node a and the point when it infects P_{crit} nodes.

Since the network structure of an unrewired lattice is isomorphic across all neighborhoods, successful contagion will spread in a uniform pattern isomorphic to the seed cluster. The seed cluster contains $k + 1$ nodes total, a seed node plus its k neighbors. For minimal threshold $\rho_{crit} = 1/k$, at each time point contagion will infect an additional k nodes. For maximum threshold that can spread on the lattice, somewhere around $1/2$ (cite - Morris), at each time point contagion will infect an additional 2 nodes. So contagion will infect P_{crit} nodes after $\rho_{crit}P_{crit}$ time steps.

At the same time, if the local information set of a is $L(a)$, contagion will infect a shortly after it infects its immediate network neighbors. The distance d between the seed and a 's network neighbors is $D(a, E(A)) - 1$. In each time step, contagion covers between $1/d$ (for $\rho_{crit} = 1/k$) and $k/(2d)$ (for $\rho_{crit} = 1/2$) of that distance: of the nodes it infects in each time step, exactly one half form an unbroken path that is between $E(A)$ and a on both endpoints. So contagion will infect a 's neighbors after $\rho_{crit}D(a, E(A))k$ time steps. So the bias for node a will be given by:

$$Bias(L(a)) = O(\rho_{crit} |D(a, E(A))k - P_{crit}|)$$

□

It is important to note that this bias value is unknown for agents a , since they don't know how close, or far, they are to $E(A)$. We have to include randomly selected agents a_g into a 's local information set to get a determinate bound on bias.

Now consider the local information set $L'(a) = a_l[q]a_g[r]$, i.e. with exactly q

elements a_i and exactly r elements a_g . We already know what happens when $r = 0$. Therefore, it is possible to calculate the improvement in bias by considering the marginal effect of increasing r in $L'(a)$. We begin by calculating the distribution in the distance between nodes a_g and $E(A)$ as a function of r . For simplicity, we hold q constant at 0.

Lemma 2.11.2. *For a Small World graph G represented as an unrewired lattice with each node having k edges, and a contagious phenomenon seeded with one node e and its network neighborhood, so $E(A) = e \cup Nbrs(e)$, and a local information set $L'(a) = a_g[r]$, no more than $1/z^2$ of the r nodes will fall outside of z standard deviations σ from the expected value $E(P)$ of distance between a_g and $E(A)$, where:*

$$E(P) = \frac{1}{2} \left(\frac{N}{k} + 1 \right)$$

$$\sigma(P) = \left[\left(\frac{N}{k} + 1 \right) \left(\frac{1}{6} \left(\frac{2N}{k} + 1 \right) - \frac{1}{4} \left(\frac{N}{k} + 1 \right) \right) \right]^{1/2}$$

Proof. Consider the case where $r = 1$. Then $L'(a)$ contains exactly one agent a_g picked uniformly at random from the total population (since the graph is degree regular, each agent has equal probability of being picked using the random strategy). As in the previous lemma, we are interested in the distance between this randomly picked agent and $E(A)$. This distance is picked from a distribution:

$$P = P(D(a_g, E(A)))$$

Note that all agents are evenly distributed into bins of k agents each that are at

some distance D from $E(A)$. The closest bin is distance 1 away from $E(A)$ and the furthest bin is distance $N/2$ away from $E(A)$. By Chebyshev's inequality(cite), no more than $1/z^2$ agents can be further than z standard deviations σ away from the mean $E(P)$. The mean is given by:

$$\begin{aligned} E(P) &= \frac{1}{N} \sum_{i=1}^{\frac{N}{k}} ki = \frac{k}{N} \sum_{i=1}^{\frac{N}{k}} i \\ &= \frac{1}{2} \left(\frac{N}{k} + 1 \right) \end{aligned}$$

The standard deviation of P is given by:

$$\begin{aligned} \sigma(P) &= \sqrt{E[(X - E(P))^2]} \\ E[(X - E(P))^2] &= \frac{k}{N} \sum_{i=1}^{\frac{N}{k}} \left(i - \frac{1}{2} \left(\frac{N}{k} + 1 \right) \right)^2 \\ &= \frac{k}{N} \left(\sum_{i=1}^{\frac{N}{k}} i^2 - \left(\frac{N}{k} + 1 \right) \sum_{i=1}^{\frac{N}{k}} i + \frac{1}{4} \left(\frac{N}{k} + 1 \right)^2 \right) \\ &= \frac{k}{N} \left(\frac{1}{6} \frac{N}{k} \left(\frac{N}{k} + 1 \right) \left(\frac{2N}{k} + 1 \right) - \frac{1}{2} \left(\frac{N}{k} + 1 \right) \left(\frac{N}{k} \right) \left(\frac{N}{k} + 1 \right) + \frac{1}{4} \left(\frac{N}{k} + 1 \right)^2 \right) \\ &= \left(\frac{N}{k} + 1 \right) \left(\frac{1}{6} \left(\frac{2N}{k} + 1 \right) - \frac{1}{4} \left(\frac{N}{k} + 1 \right) \right) \\ \sigma(P) &= \left[\left(\frac{N}{k} + 1 \right) \left(\frac{1}{6} \left(\frac{2N}{k} + 1 \right) - \frac{1}{4} \left(\frac{N}{k} + 1 \right) \right) \right]^{1/2} \end{aligned}$$

□

We now have all the tools to prove the following theorem:

Theorem 2.11.3. *For a Small World graph G represented as an unrewired lattice with each node having k edges, and a contagious phenomenon seeded with one node e and its network neighborhood, so $E(A) = e \cup Nbrs(e)$, populating the local information set $L'(a)$ of a with r agents picked at random from the population (a_g) will at best decrease the bias of a as $O(\sqrt{r})$ up to the theoretical maximum threshold $\rho_{crit} = 1/2$.*

Proof. Recall from Lemma 10.1 that when the local information set of a is $L(a)$, the bias is a function of P_{crit} and the distance between $E(A)$ and a . In the case where $r > 0$ (and q as in Lemma 10.2, set to 0), the bias is no longer a function of the distance between $E(A)$ and a since a no longer relies on its immediate network neighbors when deciding whether to adopt the contagious phenomenon. Instead, a relies on some random nodes a_g and will become infected when these nodes become infected, or not at all. When $r = 1$, the bias of node a is given by:

$$\begin{aligned} Bias(L'(a)) &= O\left(p_{crit} \left| D(a_g, E(A))k - P_{crit} \right|\right) \\ &= O\left(p_{crit} \left| E(P)k - P_{crit} \right|\right) \end{aligned}$$

The advantage of this new bias value is that it is known to the agents, because it can be calculated from known information (N , k , and P_{crit}). The agent can now estimate the impact of increasing r on the bias. As more and more randomly picked agents are added to this local information set, Chebyshev's inequality states that some very few of them will fall far outside the mean. In particular, in a set of r agents, Chebyshev's inequality states that at most one agent ($1/r$) will fall \sqrt{r} standard deviations away from the mean. Since the mean $E(P)$ and the standard deviation $\sigma(P)$ are both $O(N/k)$, then a choice of r randomly picked agents will yield at most 1 agent that lies outside the following range of

distances from $E(A)$:

$$\left(\frac{1}{2}\left(\frac{N}{k} + 1\right) - \sqrt{rc}\frac{N}{k}, \frac{1}{2}\left(\frac{N}{k} + 1\right) + \sqrt{rc}\frac{N}{k}\right)$$

for some constant c . The one-sided version of Chebyshev's inequality allows us to pick one side of that range (for instance, picking all the nodes that are more than \sqrt{r} standard deviations closer to $E(A)$ than the mean, by changing r to $r + 1$ (no significant impact on sample size).

First let's consider the case where $P_{crit} \ll E(P)k$. Then most randomly picked agents will lie relatively far away from $E(A)$, so relying on their information will prevent a from adopting early enough, keeping the bias at the expected value for $r = 1$ or even increasing it. As r goes up, however, a very small number of agents will be located closer than expected to $E(A)$. Since the threshold is very low, even one of these agents might bring a over the adoption threshold earlier than it would have by relying only on a few randomly picked agents, and thus decrease the bias. A linear increase in the size of r will bring a square-root increase in the number of these closer nodes, and a square-root decrease in the bias.

Now consider the case where $P_{crit} = E(P)k$. Then:

$$\rho_{crit} = P_{crit}/N = E(P)k/N \approx 1/2$$

which is the theoretical maximum threshold for a contagious phenomenon to spread on an unrewired lattice, as cited above.

□

We show the behavior of a 's bias as a function of P_{crit} and r for specific values of these parameters in the simulation section. We also investigate interesting behavior linked to the usage of nodes with $L'(a) = a_g^*$ as "shortcuts" in the network.

We conclude by investigating the case where $p > 0$ so some of the ties are rewired. In this case, the network is truly a small world, as it maintains some degree of clustering, but also gains high connectivity.

Theorem 2.11.4. *For a Small World graph G represented as a rewired lattice with each node having k edges and edge rewiring probability p , and a contagious phenomenon seeded with one node e and its network neighborhood, so $E(A) = e \cup Nbrs(e)$, for $\rho_{crit} \leq 1/k$, the phenomenon spreads throughout the entire network in a logarithmic number of steps, so the average number of sub-optimal time steps for nodes in $S(A)$ is minimal given original local information set $L(a)$.*

Proof. A Small World graph has the connectivity of a Poisson Random graph, and, as in a Poisson Random Graph, the targets of the rewired ties are random. Since degree remains constant throughout rewiring, at $\rho_{crit} \leq 1/k$ the contagious phenomenon can spread through any tie in the network, and will follow the spread pattern of a successful contagious phenomenon on a Poisson Random Graph, infecting all nodes in a logarithmic number of steps, as per the argument in Theorem 2.10.1 □

Finally, there is the case of Small World graphs where $1/k < \rho_{crit}$. I do not present a formal proof for this case. Instead, I describe a general outline for a 's strategies. In this case, as I will show in the next chapter, contagion goes through two phases: a "ramp-up" phase when it moves through short-range

ties outward from $E(A)$ and a “critical” phase when contagion begins to take advantage of the rewired ties. In the second phase, contagion behaves exactly as *simple* contagion with $\rho_{crit} \leq 1/k$. The crucial question then becomes: does the “critical” phase begin before, after, or exactly when P_{crit} nodes have become infected? In the first case, it is a 's goal to become infected later in the critical phase than it would normally be infected if relying on $L(a)$. In this case, a can minimize bias by increasing r : even though most nodes will become infected very quickly, merely by increasing the number of nodes in its local information set a can increase the chances that it will discover some agents that become infected later than others, and so delay its adoption time step. In the second case, it is a 's goal to become infected before the critical phase. Again, a can accomplish this by raising r , as it increases the chances of including nodes in $E(A)$ or close by in its information set. In the last case, the rest of the nodes become infected shortly after the “critical” phase begins, and we can use reasoning similar to Theorem 10.4 to show that $1 - P_{crit}$ nodes behave optimally with their original local information sets $L(a)$.

2.12 Preferential Attachment Graph

The last model graph case we consider is the preferential attachment graph [12]. This graph model differs from the previous two: unlike the random graph, ties are not created randomly between two nodes but rather in proportion to their degree. This means some redundancy exists in the graph structure and I cannot apply my analysis of local information sets on random graphs directly. At the same time, the preferential attachment graph is far from degree-regular, so even for a fixed value of ρ_{crit} across all agents, a very few high-degree agents will take

a lot of infected neighbors before they themselves become infected. As a result, I cannot apply my analysis of local information on small world graphs. Instead, I focus on a different analytical approach rooted in the existence of very-high degree agents in preferential attachment graphs and the importance of those agents for the graph's connectivity.

I begin with a simple lemma:

Lemma 2.12.1. *For a Preferential Attachment graph G with maximum agent degree k_{max} , and a contagious phenomenon with threshold $\rho_{crit} \leq 1/k_{max}$, the number of time steps any node behaves suboptimally is logarithmic in N , the number of nodes in G .*

Proof. When $\rho_{crit} \leq 1/k_{max}$, any agent in the network will become infected with having as little as one infected neighbor. In essence, the contagious phenomenon can leverage any and all ties, starting with the ties out from $E(A)$, to spread throughout the graph. Thus, the contagious phenomenon will expand outward from $E(A)$ and reach any node a in $S(A)$ in the number of time steps equal to the number of hops in the shortest path between $E(A)$ and $S(A)$. In a preferential attachment graph, the length of this path is logarithmic in G , so a will become exposed to (and adopt) the contagious phenomenon in $O(\log(N))$ time steps. \square

This lemma is helpful for setting a lower limit for the spreadability of contagion on preferential attachment graphs, but the limit it places is in most cases too low to be realistic. Given the highly skewed degree distribution of these graphs, $1/k_{max}$ can easily be as low as $1/1000$ or less. Behaviors with such low threshold may be very rare in nature. What happens when a higher-threshold contagious phenomenon spreads on a preferential attachment network? To shed light on

this question, I prove the following theorem:

Theorem 2.12.2. *For a Preferential Attachment graph G with maximum agent degree k_{max} , and a contagious phenomenon with threshold $\rho_{crit} > 1/k_{max}$, the contagious phenomenon begins by infecting nodes with lower degree to accumulate sufficient reinforcement ties to infect nodes with higher degree. Specifically, if the contagious phenomenon can infect nodes with degree less than d with just one exposure, then with likelihood at most $\Xi\phi_l(G)$ it can eventually overcome the threshold:*

$$\rho_{critmin} = \frac{\Xi}{Nk_{max}^{2-\alpha} - 1}$$

where:

$$\begin{aligned}\phi_l(G) &= 1 - (1 - f_g(d))^l \\ \Xi &< |E(A)|((1 - c)d)^l \\ f_g(d) &= O\left(\frac{1}{d^{\alpha-2}}\right)\end{aligned}$$

In the equations above, c is a clustering coefficient parameter that gives the asymptotic clustering coefficient for a node with degree d in a Preferential Attachment graph, and l is a parameter indicating the length of a path outward from $E(A)$ that contains no node with degree d or greater.

Proof. I first outline the argument for this proof. When $\rho_{crit} > 1/k_{max}$, the highest-degree agents in the network need multiple ties to infected individuals before they themselves become infected. In its initial stages, the contagious phenomenon must “sidestep” these agents as it spreads outwards from $E(A)$ and

infects lower-degree agents first. However, there is an important obstacle in the way of this further diffusion - due to the nature of the preferential attachment graph, all paths from $E(A)$ to other agents in G are most likely to lead through higher degree agents. A tradeoff emerges: the further the contagious phenomenon spreads outward from $E(A)$, the less likely it is to spread further, since all paths are “blocked” by high-degree agents, but the more agents it has infected, the more likely the contagious phenomenon is to have sufficient reinforcing ties to overcome the threshold of high-degree agents.

I now formalize this argument.

Consider a path P outward from $E(A)$ of length l . This path is much more likely to go through higher-degree nodes than lower-degree nodes. The frequency f of these nodes is given by the definition of a preferential attachment graph. Specifically, the likelihood $f(d)$ of an edge pointing to a node with degree d is:

$$f(d) = \frac{d|\{i \in G : deg(i) = d\}|}{SDG}$$

where SDG is the sum degree over all nodes in G . In the limit, we can approximate this quantity in terms of the degree distribution of G , which follows a power law:

$$\begin{aligned} f(d) &\approx \frac{dP(d)}{\int_1^{k_{max}} xP(x)dx} \\ &= \frac{d^{1-\alpha}}{\int_1^{k_{max}} x^{1-\alpha}dx} \end{aligned}$$

now consider the likelihood $f_g(d)$ of an edge pointing to a node with degree d or greater, which follows from the above:

$$\begin{aligned}
 f_g(d) &\approx \frac{\int_d^{k_{max}} x^{1-\alpha} dx}{\int_1^{k_{max}} x^{1-\alpha} dx} \\
 &= \frac{k_{max}^{2-\alpha} - d^{2-\alpha}}{k_{max}^{2-\alpha} - 1} \\
 &= O\left(\frac{1}{d^{\alpha-2}}\right)
 \end{aligned}$$

Having the baseline probability of any edge in G pointing to a node with at least degree d , we can calculate the same for any edge in a sequence of l edges. Let's call this quantity $\phi_l(G)$:

$$\phi_l(G) = 1 - (1 - f_g(d))^l$$

Going back to the original problem, $\phi_l(G)$ tells us the likelihood that a path of length l will have at least one node with degree at least d . Now consider all paths out of $E(A)$ of length l . The number of such paths is bounded by the sum degree of the nodes in these paths (at most $d - 1$) and the clustering coefficient of these nodes, which gives the fraction of paths that "loop back" on each other. Specifically, the number Π of paths of length l out of $E(A)$ is bounded by:

$$\Pi < |E(A)|((1 - c)d)^{l-1}$$

where the clustering coefficient c of a node with degree d is $O(d^{-1})$ [65]. Note that this bound is only valid when c is low, as it is for Preferential Attachment

graphs. The expected fraction of the paths that have at least one node with degree at least d is bounded by:

$$\Pi\phi_l(G) < 1 - (1 - f_g(d))^l |E(A)|((1 - d^{-1})d)^l$$

This quantity gives the expected number of paths of length l that the contagious phenomenon will spread through while avoiding nodes of degree d or greater. Similarly, we can bound the number of nodes that lie within l steps from $E(A)$ as

$$\Xi < \Pi(1 - c)d$$

We can use the same formula as above to give the fraction $t(d)$ of these that point to nodes with degree d :

$$t(d) = \frac{dP(d)}{\int_1^{k_{max}} x^{1-\alpha} dx}$$

Note that we here want to focus on nodes with degree exactly d which have the greatest chance of being infected through reinforcing ties, as opposed to all nodes with degree d or greater. The number of nodes within l distance from $E(A)$ that point to nodes with degree at least d is, in turn, given by $\Xi t(d)$. Meanwhile, the number of nodes with degree d in the entire network is given by $NP(d)$, and each of the source nodes is equally likely to point to a target node, so the final expected number of infected nodes pointing to a node with degree d is:

$$t(d)\Xi/(NP(d)) = \frac{\Xi d}{Nk_{max}^{2-\alpha} - 1}$$

This quantity divided by d , finally, gives the expected threshold that can be overcome by a contagious phenomenon as it avoids nodes of degree d or greater:

$$\rho_{critmin} = \frac{\Xi}{Nk_{max}^{2-\alpha} - 1}$$

□

2.13 Simulations

I conclude this section with a few simulation results that corroborate the analysis above and extend it to cover more complex network structure and threshold combinations. I begin with a simple plot of the size of infected cluster vs ρ_{crit} for a Poisson Random graph ($N = 1000$ nodes, $p = .01$). The plot shows that, indeed, below $\rho_{crit} = \frac{1}{Np}$ the contagious phenomenon takes over the entire network, while above it the contagious phenomenon takes over only a few nodes. Note that the transition does not happen exactly at $\frac{1}{Np}$ due to the stochastic nature of tie formation in the network.

Separately I show the average bias for the default local information set $L(a)$ over all nodes a in $S(A)$ (as defined in section 6) vs. ρ_{crit} . For $\rho_{crit} < \frac{1}{Np}$, the bias is small, on the scale of $\log(N)$. Above this value, bias is 0 as the contagious phenomenon never spreads beyond the seed nodes and so does not infect nodes in $S(A)$.

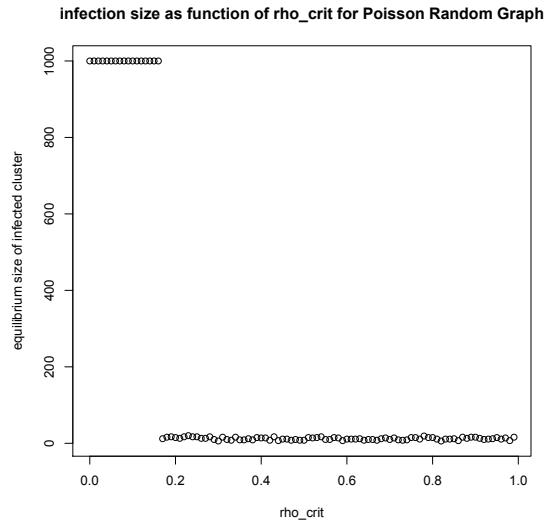


Figure 2.1: ρ_{crit} vs. size of infected cluster for Poisson Random Graph

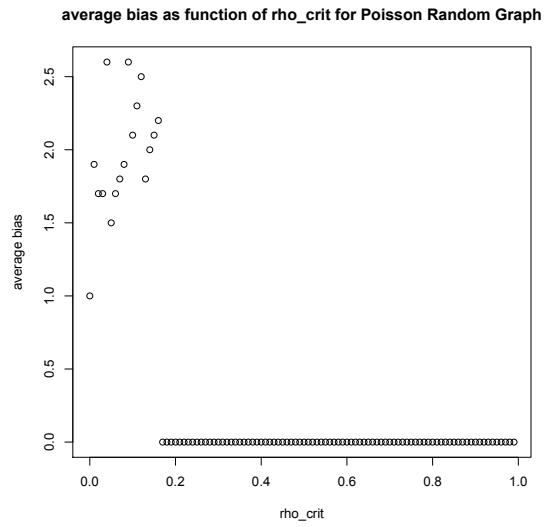


Figure 2.2: ρ_{crit} vs. average bias for Poisson Random Graph

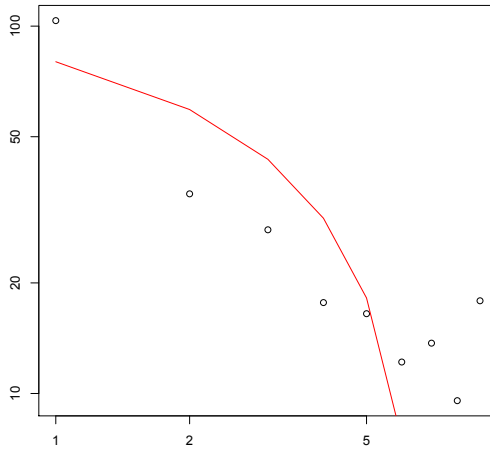


Figure 2.3: r vs. average bias for Unrewired Lattice Graph

We now examine the unrewired lattice graph ($N = 1000$ nodes, $k = 4$) and show the effect of increasing r on average bias, for $\rho_{crit} = 1/k$. The points show average bias vs. the value of r on log-log axes (points), the best fit to the data, and a fitted line of the form $y \approx \sqrt{(x)}$ (red line). The first set of results show a dramatic decrease in bias as r increases that greatly exceeds the expectations set by analysis. After further investigation, I have discovered why this is the case: in the model as written, a node in $S(A)$ can, once infected, go on to infect its neighbors if the ρ_{crit} is sufficiently low. As a result, nodes in $S(A)$ that by chance have their local information sets close to the seed nodes act as “shortcuts” in the network, accelerating the spread of the contagious phenomenon. This has two effects: one, it decreases the time to infect P_{crit} nodes. Two, it decreases the time to infect other nodes in $S(A)$. As a result, bias shrinks overall.

To correct for this effect, I changed the simulation to never spread through nodes in $S(A)$ - so nodes in $S(A)$ can become infected, but can never infect their neigh-

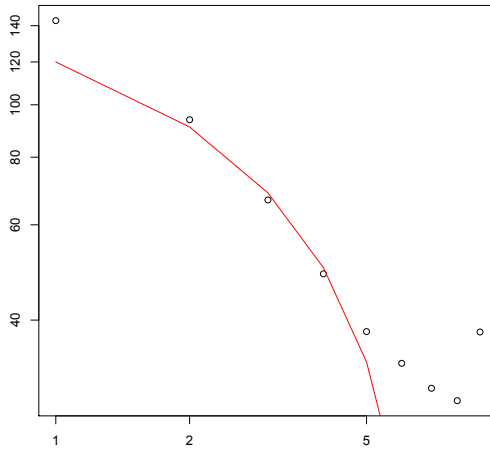


Figure 2.4: r vs. average bias for Unrewired Lattice Graph without “shortcuts”

bors. This change prevents these nodes from acting as shortcuts without dramatically changing the other aspects of the diffusion process (so long as the size of $S(A)$ is small, and ρ_{crit} low, the contagious phenomenon can essentially ignore these nodes as it spreads throughout the network). After changing the simulation, I reran it to generate a new plot of average bias vs. r as shown in Figure 2.4. The plot shows a sublinear decrease in bias as a function of r that is consistent with the square root function (again, a fitted function of the form $y \approx \sqrt{x}$ is shown as a red line). For higher values of r , the decrease becomes slower yet, indicating that while \sqrt{r} remains an upper bound for bias decrease, other factors may contribute to an even slower decrease in bias. It remains an open question what these factors are and how to account for them.

Finally, I plot the time to infect all nodes vs. ρ_{crit} for a Barabasi-Albert preferential attachment graph ($N = 1000$ nodes, $m = 5$). I vary ρ_{crit} from $1/k_{max}$ to $1/k_{avg}$, the average network degree. As expected, the time to infect all nodes t

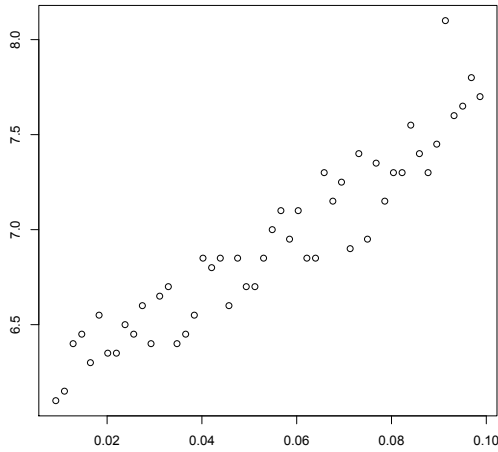


Figure 2.5: time to infect all nodes vs. ρ_{crit} for Barabasi-Albert graph

is logarithmic in N and goes up linearly as ρ_{crit} goes up due to the delay factor of accumulating sufficient reinforcement ties to infect the highest-degree nodes. As a result, average bias over this set of ρ_{crit} values remains logarithmic in N .

2.14 Discussion

The analytic and simulation results suggest several high-level observations about strategies for optimizing local information sets of agents to adopt as close as possible to the globally-optimal time step. The first observation is that network structure often plays a more important role than threshold. For instance, on a Poisson Random graph, the optimal strategy is to maintain one's original local information set $L(a)$ regardless of the threshold ρ_{crit} . In contrast, on an unrewired Small World graph, the optimal strategy is to query nodes at random from the population, again across the range of ρ_{crit} values. This observation sug-

gests that real-world applications of this model should pay attention to network structure at least as much as if not more than threshold value: trying to optimize the adoption behavior of agents in a poorly-connected network with few hubs is very different than trying to optimize the behavior of individuals in a well-connected network with many hubs, both for high-threshold and low-threshold contagion.

The second observation is that for well-connected networks with little local structure (such as the Poisson Random Graph and the Preferential Attachment Graph), the original local information set $L(a)$ is often an optimal local information set. Low-threshold contagion on these networks spread in logarithmic time, making any potential improvement from changing local information sets effectively insubstantial. In contrast, high-threshold contagion will often not spread at all, as these networks lack the redundant ties to facilitate their spread.

The third observation is that the analysis and simulations are closely linked to the somewhat stylized structure of the model. Empirical networks exhibit a complex interplay of local structure and global interconnectivity that is hard to capture in a model network. Individual decision-making often does not follow strict utility functions. Tie strength is not homogeneous in empirical networks, and people in the real world may pay more attention to the decisions of their close friends than to those of acquaintances. It is important to be aware of these differences and not assume that the results of analysis and simulation in this chapter will carry over to real-world agents. More research is required, both on individual decision-making patterns, on the structure of social networks, and on the nature of social contagion in general, before it is possible to translate models into experiments and policy implications. Still, I remain confident that the core

principles outlined in this chapter will survive, in some form, application to real-world problems of optimal decision-making, such as misinformation contagion.

CHAPTER 3

CHAPTER 3: COMPLEX CONTAGION AND CRITICAL MASS

Why do some contagion “go viral” and others do not? Research on “small world” networks [72] shows how a very small number of long-range ties that bridge between clusters can allow contagion to spread almost as rapidly as on a random network of equal density. Recent research shows how long-range ties that accelerate the spread of information and disease can impede the spread of *complex contagion*—behaviors, beliefs and preferences that diffuse via social influence and therefore often require contact with multiple adopters [21]. In confirming this result analytically and extending the analysis from small world to power law networks, we discovered that complex contagion requires a critical mass of infected nodes that corresponds to a phase transition in the ability of the contagious phenomenon to take advantage of the “shortcuts” created by long-range ties. We demonstrate how this critical mass is related to the dynamics of the contagious phenomenon and identify implications for modeling behaviors that spread via social influence, such as viral marketing and social movements.

3.1 Background

“Critical mass” refers to the point at which a dynamical process becomes self-sustaining, as in nuclear fission. Our interest focuses on critical mass in the spread of complex contagion [21] – behaviors, beliefs, and preferences that spread via social influence and therefore require contact with multiple adopters. Examples include the tipping point in the adoption of a virally marketed product, the take-off point at which a new fashion becomes a fad, and the sudden

explosion of participation in a new social movement whose appeal increases with its popularity. For complex contagion, critical mass is the level of infection at which the contagious phenomenon becomes self-sustaining, such that each additional infected node leads to one or more additional nodes to become infected, until the contagious phenomenon saturates the target population. We model the dynamics of critical mass in complex contagion and identify the critical value at which the spread of complex contagion phenomenon becomes self-sustaining. We show how this value can be derived analytically, and how the value depends on network topology and the network externality of the contagious phenomenon.

In contrast to simple contagion, complex contagion have a critical value at both the individual and the network levels. At the individual level, network externality refers to the effect of infection of some node in a network on the probability that adjacent nodes will also become infected. The spread of many social contagion phenomena—which may be risky or costly to adopt—has a positive network externality, so the probability of a node’s infection increases with the number of that node’s infected neighbors. We focus on a very simple model of positive network externality: the threshold model of contagion, where a node will become infected if and only if a critical fraction of its neighbors have become infected. This property has important implications for the population dynamics: complex contagion phenomena may die out early on due to a lack of redundant exposures to multiple infected neighbors. Thus complex contagion can have a critical value at both the individual level and the population level. At the individual level, the critical value is the proportion of neighbors who must be infected in order to infect a given individual, which we define as the “threshold.” At the population level, the critical value is the proportion of the

population who must be infected in order for the contagious phenomenon to become self-sustaining, which we define as the "critical mass." Our research investigates the relationship between individual-level thresholds and population-level critical mass in social contagion.

The population-level and individual-level critical values are known to be related [34, 62], but the form of the relationship is not well understood. We show that for two widely studied network structures – small world and power law networks – a population-level critical mass exists that may be far smaller (as a proportion of population size) than is the individual level threshold (as a proportion of neighborhood size, where a neighborhood is defined as a node and all nodes adjacent to it). This critical mass corresponds to a phase transition, from local to global propagation, once the contagious phenomenon acquires the ability to escape the region of initial infection via long-range ties that bridge across clusters. We demonstrate how this critical mass is related to the dynamics of the contagious phenomenon. Prior to reaching critical mass, complex contagion phenomena are fragile and highly dependent on the idiosyncrasies of local network structure. However, once contagion reach critical mass, they become self-sustaining and highly likely to spread throughout a connected population. In the conclusion, we discuss the theoretical implications of our analysis for explaining how and why contagion phenomena "go viral" and the practical implications for viral marketing and social movement mobilization.

Empirical research into complex contagion is a relatively new field, however, we would like to draw attention to two recent studies. Centola [20] studies the spread of health behavior in artificially structured online communities, and finds that individual adoption was more likely when participants received

social reinforcements from multiple network neighbors and that health behavior spread farther and faster across clustered-lattice networks than across corresponding random networks. Romero et al. [60] study the diffusion of “hashtags” (which include information, memes, and proxy for certain behaviors like joining a political movement) through the Twitter social network and find that hashtags related to internal Twitter memes are less likely to be adopted after multiple exposures relative to hashtags related to politics, suggesting that political hashtags may more closely resemble complex contagion than memes. Based on these and other studies, we believe that our research into the dynamics of complex contagion has applications to the spread of behaviors and products in the real world.

3.2 Model

Our research builds on and extends the model of complex contagion developed by Centola and Macy [21], which is identical to the Watts and Strogatz small-world model [72] except that it allows for individual thresholds greater than one. The model operates on a networked population of N agents, where each agent has one binary state, “infected” or “uninfected,” representing whether the individual has adopted a behavior (e.g. acquired a new technology), belief (e.g. an urban legend), or norm (e.g. smoking is uncool). The agents in this population can change their state only from uninfected to infected, in the following deterministic way: if some threshold quantity a , or fraction z/n of an uninfected agent’s neighbors are infected, its state changes to infected as well. The quantity a is called the *absolute threshold* of infection for a node, and the fraction z/n is called the *relative threshold* of infection. We first consider the case of absolute

thresholds but explore relative threshold models towards the end of our analysis. Following Centola and Macy, we define "simple contagion" as $a=1$ for absolute thresholds and $z = 1/n$ for relative thresholds. A contagious phenomenon is "complex" if $a > 1$ or $z > 1/n$. A small set of nodes A have threshold $a = 0$, corresponding to "innovators" who will buy a product or join a movement without social influence from prior adopters in their neighborhood. Centola and Macy restrict the set A to the neighborhood of a randomly chosen node. Starting with the set A , the contagious phenomenon is propagated throughout the population until no more agents can be infected (which can happen if all agents are infected or if no uninfected agent has at least a infected neighbors (or z/n)).

3.3 Analysis

Centola and Macy focus on the spread of complex contagion on a perturbed regular lattice of degree eight, with p ties that are randomly rewired ($0 \leq p \leq 1$) in pairs so as to leave the degree of each node unchanged. If $p = 0$, all ties in the lattice have range 2 since every pair of adjacent nodes has at least one neighbor in common. A small amount of random rewiring transforms the lattice into a "small world" network characterized by a few long range ties and high average clustering coefficient. That is because random rewiring on a large lattice almost always replaces ties with minimal range (a path length of two steps) with long-range ties that create "shortcuts" for the spread of a contagious phenomenon over the lattice. These shortcuts allow simple contagion to spread throughout the network far more quickly than would be possible along highly clustered short-range ties. For example, in a small-world network of N nodes, the path length between any two nodes via short-range ties is $O(N^{1/2})$ steps, whereas the

path length via long-range shortcuts is $O(\log(N))$ steps.

Centola and Macy's contribution was to show that complex contagion phenomena spread farther and faster on an unperturbed lattice ($p = 0$) than on a small-world network (e.g. $p = .1$). The spread of these phenomena on a small world network depends on the probability of contact with additional infected neighbors, given that there is contact with one infected neighbor. For nodes whose networks are highly clustered (their neighbors are also neighbors of one another), if a node has one infected neighbor then that node is highly likely to have other infected neighbors, even when the proportion of the population that is infected is still very small. In contrast, for nodes whose neighbors are randomly chosen, when the size of the infected population is small, so too is the probability that a node will have additional infected neighbors, given that the node has one infected neighbor. As the number of infected nodes grows, so too does the probability that random ties will connect an uninfected node with a sufficient number of infected neighbors for that node to also become infected. We identify the critical mass as the point at which the proportion of infected nodes is sufficient for contagion to take advantage of the shortcuts created by long-range ties. For simple contagion, the critical mass is uninteresting, since it is achieved at a single infected node, i.e. the seed node of the contagious phenomenon. For complex contagion, the critical mass is always greater than one, and knowing how much greater is important from both a theoretical and practical perspective.

3.3.1 Small World Networks

We begin by deriving the critical mass for small world networks, modeled as a perturbed regular lattice. This model captures two defining properties of empirically observed small world networks – that most ties have minimal range (the modal range is 2 in most empirical networks), while the network also has a relatively small number of long-range ties. The short-range ties correspond to the high level of clustering that is observed in most empirical social networks, while a few long-range ties make possible the surprisingly short mean geodesic (such as the widely observed “six degrees of separation”). Rewired ties on the perturbed lattice are a highly stylized representation of the empirical regularity that nodes in a social network have some mix of highly clustered ties (usually to close friends and family) and long-range ties (usually to acquaintances). Following Centola and Macy, we impose the conservative simplification that infected acquaintances exert as much influence as close friends. Relaxing that assumption is equivalent to increasing a , the threshold number of infected neighbors that are required for a node to become infected. The derivation of the critical mass is based on the probability P_{RW_a} of an uninfected node having a rewired ties to infected nodes, as the number of infected nodes increases.

Theorem 3.3.1. *Given a randomly rewired lattice of N nodes, where every node has probability p of having one of its ties rewired, and I infected nodes on that network, the probability that any uninfected node c has a rewired ties to infected neighbors is given by:*

$$P_{RW_a} \approx 1 - P_{NIa}(c) \binom{k}{a}^{(N-I)} \quad (3.1)$$

where

$$P_{NIa}(c) = 1 - p^a + p^a \left(\frac{\binom{N-1}{a} - \binom{I}{a}}{\binom{N-1}{a}} \right) \quad (3.2)$$

Proof. We can interpret P_{RW_a} as one minus the probability that no node has a rewired ties to infected nodes:

$$P_{RW_a} = 1 - \prod_{(c) \in V \setminus F} P_{NI}(c) \quad (3.3)$$

where V is the set of all nodes and F is the set of infected nodes. Then, $P_{NI}(c)$ is the probability that c does not have rewired ties to a infected nodes. This probability is approximately uniform over all c on a randomly rewired lattice, except in the case where $N - I$ is very small so the number of possible targets that are not infected nodes is quickly exhausted. Given this qualification, we can rewrite the above as:

$$P_{RW_a} \approx 1 - P_{NI}(c)^{N-I} \quad (3.4)$$

$NI(c)$ holds if we can't pick a of c 's ties such that all a are rewired and both point to infected nodes. There are $\binom{k}{a}$ independent ways to pick a of c 's ties, where k is c 's degree (uniform over all nodes in a lattice). So we can again rewrite:

$$P_{RW_a} \approx 1 - P_{NIa}(c)^{\binom{k}{a}(N-I)} \quad (3.5)$$

where $P_{NIa}(c)$ is the probability that, having picked some set of c 's ties with a elements, at least one of these ties is not rewired and/or does not point to an

infected node. Tie rewiring is an independent process (the probability of one tie being rewired does not affect the probability of other ties being rewired), so with probability $1 - p^a$ at least one tie is not rewired and $NIa(c)$ holds. In the opposite case, $NIa(c)$ still holds so long as the targets of all a ties are not in F . Note that the number of possible targets of c 's ties is $N - 1$, since c can't have ties to itself. More formally:

$$P_{NIa}(c) = 1 - p^a + p^a \left(\frac{\binom{N-1}{a} - \binom{I}{2}}{\binom{N-1}{a}} \right) \quad (3.6)$$

□

For small a , we can approximate $P_{NIa}(c)$ as follows:

$$P_{NIa}(c) \approx 1 - \left(\frac{pI}{N} \right)^a \quad (3.7)$$

This theorem calculates the most conservative case where a contagious phenomenon has to spread to some node entirely via long-range ties, even though it is possible for the contagious phenomenon to spread through any combination of long-range (rewired) and short-range (unrewired) ties. Consequently, the results in this paper slightly understate the point at which a contagious phenomenon begins to take advantage of shortcuts in the network. We follow the conservative approach for two reasons: greater simplicity of analysis, and (more importantly) the implication of spreading entirely through long-range ties for the growth rate of the contagious phenomenon, which we discuss towards the end of the paper.

3.3.2 Power Law Networks

Theorem 3.1 can be extended beyond the rewired lattice. We now present a lemma that derives an approximation to P_{RW_a} for power law networks. In power-law networks, the notion of a “rewired” tie does not apply. Instead, following the Barabasi-Albert model of power-law networks [12], we assume that ties are formed according to preferential attachment, with higher-degree nodes more likely to be the targets of ties.

Lemma 3.3.2. *Given a power-law network of N nodes with sum degree NI where degree follows a power-law distribution and ties are formed according to preferential attachment, and I infected nodes with sum degree SI on that network, the probability that any uninfected node c has a ties to infected neighbors is approximated by:*

$$P_{LR_a} \approx 1 - \left[1 - \left(\frac{SI}{SN} \right)^a \right]^{k\alpha^{NTa}} \quad (3.8)$$

where k is a factor parameter given by:

$$k = \frac{a^a}{a!} \quad (3.9)$$

and NT is a parameter estimated from the degree distribution by:

$$NT \approx (N - I) \frac{a^{-\alpha+1}}{\text{mindeg}} \quad (3.10)$$

with mindeg the minimum degree in the network and α the power law exponent.

Proof. The proof is very similar to the proof of Theorem 3.1, so here we con-

centrate only on the differences. First, $P_{NIa}(c)$ no longer depends on rewiring. Without preferential attachment, we could write $P_{NIa}(c)$ as simply:

$$\begin{aligned}
 P_{NIa}(c) &= \frac{\binom{N-1}{a} - \binom{I}{a}}{\binom{N-1}{a}} \\
 &= 1 - \frac{\binom{I}{a}}{\binom{N-1}{a}} \\
 &\approx 1 - \left(\frac{I}{N}\right)^a
 \end{aligned}$$

The only correction we have to make is related to the power-law degree distribution, where each node gets up-weighted by its degree. We can transform these weights into discrete values by counting each node as many times as it has degree. Combinations of nodes from the resulting augmented sets are equivalent to weighted combinations from the original sets. The resulting equation is almost the same, except we substitute SI and SN , the sum degrees of all I infected nodes and all N nodes in the full population:

$$P_{NIa}(c) \approx 1 - \left(\frac{SI}{SN}\right)^a \quad (3.11)$$

The other difference from the proof of Theorem 3.1 is that $P_{NI}(c)$ is no longer uniform over all nodes, as they do not all have the same degree. This results in the upper exponent $N - I$ being replaced with a sum SE of $N - I$ terms, where each term represents all the ways to choose a nodes from all the neighbors of a particular uninfected node:

$$SE = \sum_{i \in V \setminus F} \binom{N(i)}{a} \quad (3.12)$$

Here $N(i)$ is the number of network neighbors of a particular uninfected node i . We now examine the terms in this sum. Each of these terms is a fraction:

$$\binom{N(i)}{a} = \frac{N(i)!}{a!(N(i)-a)!} \quad (3.13)$$

we can take $1/a!$ out of the sum, and transform as follows:

$$SE = \frac{1}{a!} \sum_{i \in V \setminus F} N(i)(N(i)-1)\dots(N(i)-a+1) \quad (3.14)$$

for small a and large $N(i)$ (which will dominate the sum), we can approximate as follows:

$$SE \approx \frac{1}{a!} \sum_{i \in V \setminus F} (N(i))^a \quad (3.15)$$

Note that we can ignore all terms in this sum where i has fewer than a neighbors (since there are no ways to choose a units from a set smaller than a , those terms are 0). Since the degree distribution follows a power-law, the number of uninfected nodes with degree $\geq a$ is given by:

$$NT \approx (N-I) \frac{a^{-\alpha+1}}{\text{mindeg}} \quad (3.16)$$

where mindeg is the minimum degree for any node in the network (we can crudely estimate it as 1) and α is the exponent of the power law distribution.

So there are NT terms in the sum overall. In a power-law distribution with discrete values, term density thins out at a rate proportional to the exponent, that is, individual degree values will be roughly powers of the exponent a . The smallest of the relevant values lies somewhere between a and aa . Since a will be by far the smallest term in the sum, we can drop it and approximate as follows:

$$SE \approx \frac{1}{a!} \sum_{j=1}^{NT} a^a \alpha^{aj} \quad (3.17)$$

Finally, we can extract a^a and the sum becomes a geometric series:

$$\begin{aligned} SE &\approx \frac{a^a \alpha^{a(NT+1)} - 1}{a! (\alpha^a - 1)} \\ &\approx \frac{a^a}{a!} \alpha^{NTa} \end{aligned}$$

□

3.3.3 Critical Behavior

Number of Infected Nodes

We now examine the behavior of P_{RW_a} at limiting values of I . First, let us consider $I < a$, where a threshold a contagious phenomenon cannot spread. For the rewired lattice we have:

$$P_{NIa}(c) = 1 - p^a + p^a \frac{\binom{N-1}{a} - 0}{\binom{N-1}{a}} = 1 \quad (3.18)$$

and for the power-law network we have:

$$P_{NIa}(c) = 1 - \frac{0}{\binom{N-1}{a}} = 1 \quad (3.19)$$

so $P_{RW_a} = P_{LR_a} = 0$ for all other parameter values. So the contagious phenomenon is indeed guaranteed not to spread through long-range ties when too few nodes are infected, because there are too few infected nodes to who uninfected nodes might be tied. In the opposite case, when $I = N - 1$, for rewired lattice we have:

$$P_{RW_a} = 1 - (1 - p^a)^{\binom{k}{a}} \quad (3.20)$$

which indicates that the probability of the final node being infected depends only on that node having a rewired ties, as the model suggests. For the power-law case, $P_{NIa}(c) = 0$, since it no longer depends on rewiring, hence:

$$P_{LR_a} = 1 - 0^{SE} \quad (3.21)$$

If $SE > 0$ (the last uninfected node has degree a or more), then $P_{RW_a} = 1$, since the node is guaranteed to have a infected neighbors. If $SE = 0$, then $P_{RW_a} = 1 - 0^0 = 0$, since the node has insufficient ties to become infected.

Threshold

Next, we consider the behavior of P_{RW_a} at limiting values of a for both the small world (perturbed lattice) and power law networks. For the case $a = 0$ for rewired lattices we have:

$$P_{NIa}(c) = 1 - p^0 + p^0 \frac{1-1}{1} = 0 \quad (3.22)$$

and for the power-law case we have:

$$P_{NIa}(c) = \frac{1-1}{1} = 0 \quad (3.23)$$

So, assuming $p > 0$ (for the rewired lattice) and $I < N - 1$, $P_{RW_a} = P_{LR_a} = 1$, which shows that a contagious phenomenon with threshold 0 is guaranteed to spread on all networks.

Now consider the case of simple contagion with threshold $a = 1$. For rewired lattices we have:

$$P_{RW_a} = 1 - \left(1 - \frac{pI}{N-1}\right)^{k*(N-I)} \quad (3.24)$$

and for the power-law case we have:

$$P_{RW_a} \approx 1 - \left[1 - \frac{SI}{SN}\right]^{\alpha^{N-I}} \quad (3.25)$$

This reduction indicates that the spread of simple contagion via long-range ties is unproblematic. For rewired lattices, simple contagion will spread across any

rewired ties between infected and uninfected nodes, even if there is only a single infected node in the population. Moreover, for an infected cluster of a given size, a simple contagion phenomenon on a power-law network is more likely to spread through long-range ties than the same phenomenon on a rewired lattice, due to the presence of nodes with very large degree in the power law distribution.

Rewiring Probability

It is also instructive to consider values of P_{RW_a} for rewired lattices for limiting values of p . For $p = 0$, we have $P_{NI_a}(c) = 1$, so $P_{RW_a} = 0$ and the contagious phenomenon cannot take advantage of long-range ties, because there are none. For $p = 1$, we have:

$$P_{NI_a}(c) = 1 - \frac{\binom{I}{a}}{\binom{N-1}{a}} \quad (3.26)$$

So the spread of the contagious phenomenon depends entirely on the size of the infected cluster.

3.3.4 Estimation of Function Behavior

We analyze other critical points of P_{RW_a} and P_{LR_a} through numerical estimation rather than by calculating precise solutions. The functional forms for these two probability functions are not readily analyzable, but numeric estimates of the functions show a number of interesting properties.

Figure 3.1 shows P_{RW_a} as a function of the number of infected nodes I for a particular set of parameters, $k = 48$, $a = 2$, $p = .1$, $N = 40000$. There are two important features to note. First, there is an inflection point in P_{RW_a} as it goes from ≈ 0 to ≈ 1 . This inflection point happens early in the contagious phenomenon diffusion process, with between 10 (or .025%) and 100 (or .25%) nodes infected. Second, there is a rapid drop-off in P_{RW_a} for very high values of I , when almost all nodes are infected. Between the inflection point and the drop-off, the value of P_{RW_a} is very close to 1.

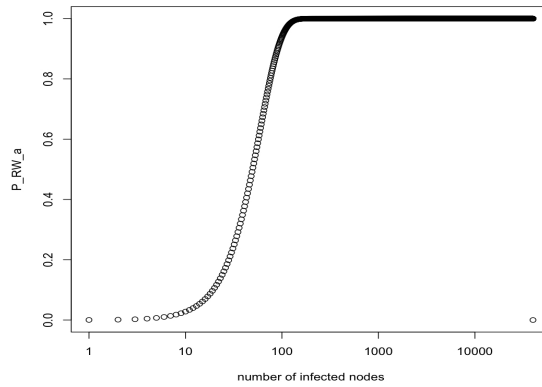


Figure 3.1: Inflection and drop-off points in the probability of **a** rewired ties to an infected node as the level of infection increases on a rewired lattice, **a = 2**, **k = 48**, **p = 0.1**, **N = 40000**

We find that the pattern evident in Figure 3.1 persists across the ranges of p , a , and N . Below, we will explore the dynamics of P_{RW_a} for these parameters, but here we focus on the only parameter that changes in the above plot, which is I . This parameter is present in P_{RW_a} in two places: $P_{NIa}(c)$ and the exponent $N - I$. As I increases, $P_{NIa}(c)$ decreases, causing an increase in P_{RW_a} , but the exponent $N - I$ decreases as well, causing a decrease in P_{RW_a} . Therefore, at the inflection point the change to $P_{NIa}(c)$ outweighs the change to the exponent, while at the drop-off near $I \approx N$, the reverse happens. In other words, at the inflection point

the likelihood of some uninfected node having a ties to infected nodes becomes so high that the smaller pool of uninfected nodes does not bring it down. At the drop-off, the pool of uninfected nodes becomes so small that the very high likelihood of some uninfected node having a ties to infected nodes does not bring the value of P_{RW_a} up.

The results of our numerical estimation suggest that the infection process is self-sustaining between the inflection point and the drop-off. In this region, each additional infected node adds more long-range ties between infected and uninfected nodes, and makes further adoption via long-range ties more likely. The beginning of this region corresponds to a phase transition where the contagious phenomenon goes from spreading exclusively via short-range ties (because P_{RW_a} is near 0) to spreading via both long- and short-range ties (because P_{RW_a} is near 1). This analysis suggests that I^* , the value of I at the inflection point in Figure 1, corresponds to a critical mass in the size of the infected population, above which complex contagion can leverage long-range ties with a sufficiently high probability for propagation via long-range ties to become self-sustaining (limited only by the declining pool of nodes that remain uninfected).

We now focus on the critical mass phenomenon and explore its values for rewired lattice networks for a parameter space of different thresholds a and rewiring levels p . Figure 3.2 below is a heat map that shows critical mass values for a range of values of a (x axis) and p (y axis), holding k constant at 48 and N constant at 40000. Colors of the contour plot correspond to values of the critical mass: red colors indicate low values, yellow values indicate intermediate values, white colors indicate high values.

Figure 2 shows how, holding p constant, low thresholds yield a smaller crit-

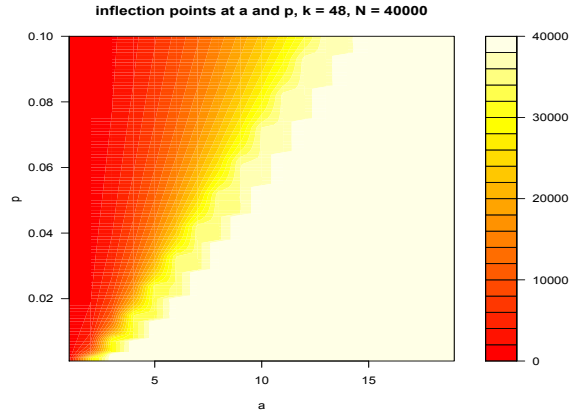


Figure 3.2: Critical mass increases with threshold a and decreases with perturbation p on a rewired lattice, $k = 48$, $N = 40000$. Colors indicate critical mass from red ($CM = 1$) to white ($CM = N$)

ical mass, and so a faster rate of spread, than high thresholds. For any $p > 0$, as a increases, the probability for an uninfected node to have a rewired ties to infected nodes necessarily decreases. For simple contagion ($a = 1$), the critical mass has its minimal value (1 infected node) regardless of p . Conversely, for contagion with very high threshold ($a \geq 15$), the critical mass has its maximum value N regardless of p . For intermediate thresholds, the critical mass decreases in p . Intuitively, the more ties that are rewired, the higher the probability that a node will have a rewired ties to infected nodes, hence a given expected value requires fewer infected nodes.

We conclude this section by replicating our analysis of P_{RW_a} on P_{LR_a} . As power-law networks have no rewiring, we add a new parameter r to represent the ratio SI/SN , which models the degree of the infected nodes relative to that of the rest of the population. Formally:

$$r = \frac{SI/NI}{I/N} \tag{3.27}$$

Figure 3.3 shows the results for a particular combination of parameter settings, $a = 2$, $N = 40000$, $\alpha = 2$ and $SI/SN = I/N$ (infected nodes have the same average degree as all nodes). We find the same overall pattern for P_{LR_a} as for P_{RW_a} with an even sharper transition (which appears as a step function due to floating-point precision limitations).

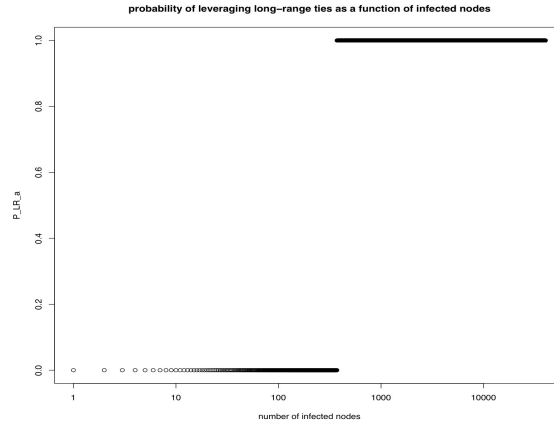


Figure 3.3: Inflection and drop-off points in the probability of **a** rewired ties to an infected node as the level of infection increases on a power law network, $\mathbf{a = 2}$, $\mathbf{r = 1}$, $\mathbf{N = 40000}$, $\mathbf{\alpha = 2}$

In Figure 3.4, we explore values of the critical mass (defined in the same way as for P_{RW_a}) for a range of values of a (x axis) and r (y axis), keeping N constant at 40000 and α constant at 2. Figure 4 shows that, for a given threshold, complex contagion phenomena on power law networks have a much smaller critical mass than complex contagion on rewired lattices. That is because the higher variance in degree in a power law network, relative to a small-world network, even with the same mean degree ($k = 8$) makes it more likely that an uninfected node will have a ties to a few hubs that have become infected.

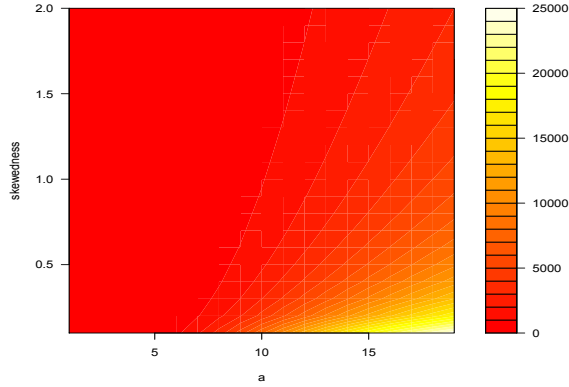


Figure 3.4: Critical mass increases with threshold a and decreases with skewedness r on a power law network, $N = 40000$, $\alpha = 2$. Colors indicate critical mass from red ($CM = 1$) to white ($CM = N$)

3.3.5 Absolute and Relative Thresholds

The preceding analysis assumed absolute thresholds. In this section, we briefly consider relative thresholds z/n . In the relative threshold model, a fraction z/n of a node's neighbors must become infected for the node to switch its state to infected. For degree-regular networks such as the rewired lattice, there is no difference between absolute thresholds a and relative thresholds z/n with respect to analysis of contagion dynamics. For networks with a non-uniform degree distribution, such as the power-law network, using a relative threshold makes the analysis more complex, but it is still possible to make broad observations about the dynamics of the contagious phenomenon. First, we have to rewrite P_{LR_a} to incorporate the relative threshold z .

Lemma 3.3.3. *Given a power-law network of N nodes with sum degree NI where degree follows a power-law distribution and ties are formed according to preferential attachment, and I infected nodes with sum degree SI on that network, the probability that any uninfected node c has z/n (where n is the node's degree) ties to infected neighbors is*

approximated by:

$$P_{LR_z} \approx 1 - \prod_{n \in \text{deg}(N)} \left[1 - \left(\frac{SI}{SN} \right)^{\lceil zn \rceil} \binom{n}{\lceil zn \rceil} n^{-\alpha} \right] \quad (3.28)$$

where $\text{deg}(N)$ is the set of all distinct degree values in the network and α the power law exponent.

Proof. The proof builds upon the derivation of P_{LR_d} but with the understanding that $P_{NI_d}(c)$ is now replaced with $P_{NI_z}(n)$ which is only uniform for nodes with the same degree. Substantively, $P_{NI_z}(n)$ is the probability that no uninfected node with degree n has zn connections to infected nodes. Note that the number of nodes with degree n in a power-law distribution is proportional to $n^{-\alpha}$.

□

To analyze P_{LR_z} , consider the relationship between $P_{NI_z}(n)$ and n . This probability has three terms, which govern its dynamics. The first is $\binom{n}{\lceil zn \rceil}$, which follows the inequality:

$$z^{\lceil zn \rceil} \leq \binom{n}{\lceil zn \rceil} \leq (ez)^{\lceil zn \rceil} \quad (3.29)$$

By the inequality, the first term is exponentially increasing in n . The second term is $n^{-\alpha}$, which is polynomially decreasing in n . The third term is $1 - \left(\frac{SI}{SN} \right)^{\lceil zn \rceil}$, which approaches 1 at an exponential rate in n . On balance, these terms indicate that $P_{NI_z}(n)$ will approach 1 at a polynomial rate in n . Furthermore, as we discussed in section 3.2, individual degree values will be roughly powers of the exponent α . These two factors in combination suggest that elements of $P_{NI_z}(n)$ for large n

will be very close to 1. That is, the very small number of nodes with high degree and the very large number of neighbors required to infect them will ensure that hubs of the network are nearly immune to infection.

It follows that the success or failure of contagion with relative thresholds depends decisively on the infection of low-degree nodes. There are two factors at play: the choice of seed cluster (hub vs. non-hub) and the interconnectivity between low-degree nodes. Even if one of the seed nodes is a hub, and the network contains multiple hubs, the contagious phenomenon can only reach these uninfected hubs by spreading through low degree nodes that are sufficiently clustered to propagate a complex contagion. This observation parallels a well-known result by Morris [49] that behaviors with a high relative threshold will spread best through local neighborhoods with a high degree of clustering. If the contagious phenomenon seed does include a hub, some degree of interconnectivity between the low-degree nodes is still necessary, assuming that the network contains multiple hubs. The fewer hubs in the seed cluster, the greater the clustering that is needed among the low-degree nodes. A single infected hub may put its low-degree neighbors over the threshold, but unless those neighbors can infect further nodes, the contagious phenomenon will not spread. We leave a more detailed analysis of relative-threshold contagion to future investigation.

3.4 Thresholds and Contagion Dynamics

We turn now from the identification of critical mass to the consequences for the propagation dynamics of complex contagion, focusing on the contagious phenomenon growth rate, that is, the proportional increase in the number of in-

fectured nodes as the contagious phenomenon spreads throughout the network. More precisely, we define the perimeter of an infected region as the number of nodes about to be infected given I nodes already infected with the contagious phenomenon, the area as the total number of nodes already infected by the contagious phenomenon, and the growth rate as the size of the perimeter relative to the area.

Definition 3.4.1. Given a contagious phenomenon with threshold a (or z/n) contagion and I infected nodes, the *perimeter* $x(I)$ of that infected set is the number of uninfected nodes that have a or more (or z/n or more) infected neighbors.

The growth rate of a contagious phenomenon over time is the ratio of the number of nodes about to be infected (the perimeter) to the number of nodes already infected, expressed as a function of I .

Definition 3.4.2. The growth rate of a contagious phenomenon $\lambda(I)$ is given by $\frac{x(I)}{I}$ as a function of I .

Consider the growth rates of a complex contagion phenomenon on a perturbed lattice before and after it reaches critical mass, starting with a single infected neighborhood A . Since the lattice has uniform degree, the analysis applies equally to absolute and relative thresholds. Before critical mass, the contagious phenomenon is unlikely to spread through long range ties since it is unlikely that an uninfected node will have a random ties to infected nodes. So, the contagion will spread by leveraging the local neighbors of A . Spatially, the neighborhood A resembles a square (for a 2-dimensional lattice; this analysis extends into lattices of dimension d). The local neighbors of A (i.e. those adjacent to A) form a perimeter around this square. As these neighbors in turn become infected, the size of the infected square increases, as does the perimeter around it.

In general, until the contagious phenomenon reaches critical mass, the infected set will always form a square with area I^2 and its perimeter will always form a square perimeter of size $O(I)$. Then the growth rate $\lambda(I)$ is given by

$$\lambda_{pre}(I) \approx \frac{I}{I^2} \approx I^{-1} \quad (3.30)$$

Thus the growth rate prior to critical mass drops quickly as I increases, for the simple reason that the perimeter of a square becomes smaller relative to its area as the area increases.

After reaching critical mass, the picture is dramatically different. The perimeter of the contagious phenomenon is not limited by the lattice structure, but also includes the expected number of nodes infected via rewired ties. Spatially, the infected set now consists of the infected square containing A plus the set of randomly distributed nodes with a random ties to infected nodes. As the size of the infected square continues to grow, one or more of these randomly infected nodes may eventually be able to help infect one of its local neighbors, and a second infected square will emerge which grows locally, as does the perimeter around it. And then a third square, and so on, each happening more quickly than the last, given the increasing overall number of infected nodes and thus the increasing probability that an uninfected node will have random ties to a infected neighbors. More formally, consider the quantity $1 - PNI(c)$, the probability of some particular node c having a or more infected neighbors:

$$1 - PNI(c) \approx 1 - \left[1 - \left(\frac{pI}{N} \right)^a \right]^{\binom{k}{a}} \quad (3.31)$$

We see that $1 - PNI(c)$ is a monotonically increasing polynomial function of I .

Now consider the quantity NIc , the expected number of nodes that have a or more infected neighbors. By an argument similar to that stated in Theorem 3.1, NIc is not exactly uniform over all c as it tends to 0 as the set of available targets (infected neighbors) is exhausted. The expected number of nodes with a or more infected neighbors will, then, be a sum of terms dominated by $1 - P_{NI}(c)$. The total number of such terms will be $N - I$, the number of nodes remaining uninfected. We can now write the post-critical perimeter as follows:

$$x(I) = O(f(I^a)(N - I)) \quad (3.32)$$

where $f(I^a)$ is some monotonically increasing polynomial function of I (i.e. $1 - P_{NI}(c)$). Then, the growth rate is given by:

$$\lambda(I) = O(f(I^{a-1})(N - I)) \quad (3.33)$$

This is a product of two functions, one monotonically increasing in I , the other monotonically decreasing in I . At first, $\lambda(I)$ is dominated by the first term (many uninfected nodes available to infect through rewired ties), and grows in I . As the number of uninfected nodes declines, the second term begins to dominate and $\lambda(I)$ falls in I . In summary, the growth rate of a complex contagion on a rewired lattice at first drops quickly in I but then, if the contagious phenomenon reaches critical mass, the growth rate suddenly takes off and increases in I and then once again drops in I as the contagious phenomenon runs out of nodes to infect.

We demonstrate this pattern empirically by plotting the observed growth rate (renormalized as $\frac{\lambda(I)}{I} + 1$ so the minimal growth rate is 1) against I for a simulated contagious phenomenon with $k = 8$, $a = 2$, $p = .1$, $N = 40000$ (Figure

3.5). The growth rate, shown in black on Figure 3.5, goes through three phases – first a rapid drop, then a sharp rise, followed by another drop. Figure 3.5 also shows (in red) the corresponding quantity $1 - P_{NI}(c) * (N - I)$, which shows the change in the probability of infection through long-range ties. Note that the latter starts growing more slowly and falls off later than the observed growth rate. This is because propagation through rewired ties alone is a conservative estimate of contagion growth and does not include growth through short-range ties, which continues after critical mass is reached even though it plays an ever-diminishing role.

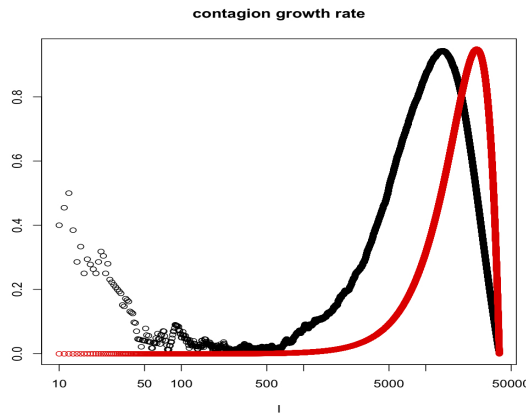


Figure 3.5: Contagion growth rate (black) and probability that a random node will be uninfected and have a random ties to infected nodes (red), for a regular lattice with $k = 8$, $a = 2$, $N = 40000$, $p = .1$

3.5 Beyond the Threshold Model

Following Centola and Macy, the preceding analyses assume discrete and deterministic thresholds of adoption. However, empirical contagion may be more plausibly modeled as continuous stochastic decisions rather than thresholds.

Previous studies [11, 41, 68] show that having multiple adopter friends does increase the likelihood of adoption, sometimes in a non-linear way. We can formalize this relationship as follows:

$$P(\text{adopt}|a \text{ adopter friends}) = f(a) \quad (3.34)$$

where f is a monotonically increasing function. For small a (region around $a = 2$ in [11]), f is convex. For large a , f is concave, with each additional adopter friend contributing a diminishing marginal likelihood of adoption.

We can use this formalism to adapt the analytical results above to contagion phenomena that do not have a deterministic threshold, but do have a positive relationship between likelihood of adoption and number of adopter friends (that is, the threshold is stochastic rather than deterministic). As in the previous section, we focus on the quantity NIc , the expected number of uninfected nodes that have a or more infected neighbors. We can simulate the behavior of NIc over values of I and a . Figure 3.6 shows the log of NIc (y axis) as a function of the log of I (x axis) on a rewired lattice with $k = 48$, $p = .1$ and $N = 40000$ and $a = 2$ (green) and $a = 3$ (brown).

Figure 3.6 shows that for values I below 10000 (1/4 of the nodes infected), NIc is an exponential function that appears linear on a logged axis. We can also see that the difference between the number of nodes that have 2 or more infected neighbors, and 3 or more infected neighbors, diminishes exponentially and disappears for $I > 10000$. This general pattern applies across the range of N , k and p , but for smaller values of k (holding all other parameters constant), the lines for $a = 2$ and $a = 3$ do not converge before the pool of uninfected nodes is

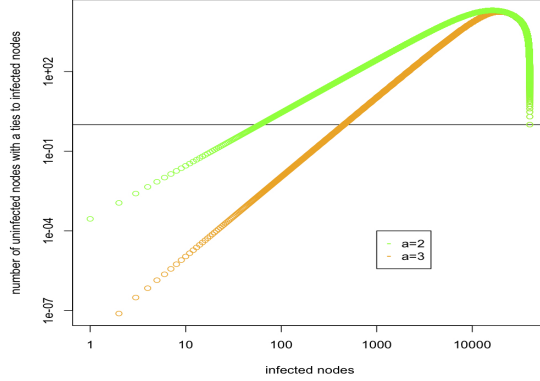


Figure 3.6: N_{Ic} (logged) as a function of the log of the area on a perturbed lattice, $k = 48$, $p = .1$, $N = 40000$, and $a = 2$ (green) and $a = 3$ (brown)

exhausted.

The dynamics of N_{Ic} as shown in Figure 3.6 suggest an important property of the diffusion of stochastic threshold contagion across regular networks: as a contagious phenomenon leverages long-ranged ties, the number of nodes that are exposed to a and not $a + 1$ infected neighbors shrinks exponentially. Given the relationship between a and likelihood of infection given by f above, this means that for small values of a , the likelihood of infection increases rapidly in the early stages of contagion diffusion. Larger values of a yield a smaller marginal likelihood, and do not need to be considered as closely. However, predictive models of contagion adoption that take into account f should outperform models that ignore the marginal likelihood.

As an illustrative example, we consider the same network as above (rewired lattice, $p = .1$, $k = 48$, $N = 40000$), and calculate the expected number of adopters at the next step for a given number of current adopters based on two different functions f . The first function, f_1 , will be linear in a , the second function, f_2 ,

will be non-linear in a . For simplicity, we consider only thresholds up to $a = 3$. The specific forms of f_1 and f_2 are as follows:

$$f_1 = ((1, .02), (2, .03), (3, .04))$$

$$f_2 = ((1, .02), (2, .05), (3, .07))$$

Then we can calculate the expected number of adopters by determining the number of uninfected nodes that have a neighbors and multiplying by the appropriate value of $f(a)$. We also include a “baseline” expected number of adopters that is based solely on $f(1)$, ignoring the marginal likelihood of adoption due to multiple exposures. The results are summarized in Figure 7:

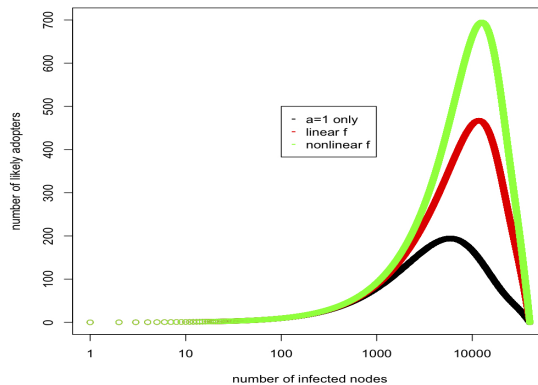


Figure 3.7: Expected number of adopters as a function of I using only $f_1(1)$ (black), all of f_1 (red), and f_2 (green), $N = 40000$, $k = 48$, $p = .1$

For $I < 1000$ and $I \approx N$, the expected number of adopters is the same whether using $f_1(1)$, all of f_1 , or f_2 . In the intermediate range $I \in [1000, N)$ (2.5% and 25% of all nodes infected), there is a significant difference in new adopter expectation. At the peak, around $I = 10000$ nodes, the “baseline” expectation is

half as much as the expectation that uses a linear f , and a third as much as the expectation that uses a non-linear f .

3.6 Discussion and Conclusion

The key contribution of this chapter is the demonstration of a bifurcation point in the spread of complex contagion – the critical mass. For simple contagion like information and disease, this bifurcation does not exist. Such contagion can leverage long-range ties even with only one infected node, hence a single seed is sufficient to create a critical mass.

For complex contagion, in contrast, the growth process will have two phases separated by a very sharp transition. Initially, the contagious phenomenon can only spread locally, that is, via short-range ties. Once every node that is reachable via short-range ties is infected, propagation terminates if the level of infection remains sub-critical.

However, if the region reachable via local propagation is sufficiently large, the contagious phenomenon will reach critical mass and the contagious phenomenon will “go viral.” On this side of the bifurcation point, the contagious phenomenon can now spread via long-range ties that allow the contagious phenomenon to “jump” to a fresh area. This area of infection then rapidly expands through short-range ties. The increase in the size of the infected population in turn increases the probability that the contagious phenomenon can spread to yet another fresh region of the network via long-range ties. Simply put, once a complex contagion phenomenon reaches critical mass, it begins to spread in the same way as a simple contagion phenomenon – taking advantage of shortcuts

to distant regions and eventually reaching every node in a connected network.

Our analysis also has an important theoretical implication for understanding why some contagion phenomena “go viral” and others do not. For a simple contagion phenomenon to escape the region of initial infection, it need only reach a node with a long-range tie or a hub that can broadcast the contagious phenomenon more widely. On a small world or undirected power law network, that is guaranteed to eventually occur, so long as the contagious phenomenon remains capable of passing from one node to another (e.g., there is no decline in infectiousness such as might happen in a news cycle). For a complex contagion phenomenon to go viral, it must infect sufficient nodes that a long-range tie or susceptible hub can make a difference.

The existence of a bifurcation point in the propagation of complex contagion phenomena has a potentially valuable practical implication for the ability to predict the eventual outcome at the early stages of a viral marketing campaign. Prior to critical mass, the growth rate decays, but as soon as the contagious phenomenon is able to spread via long-range ties, the growth rate reverses and rapidly accelerates. This qualitative change in the rate of growth from negative to positive is a statistical signature of critical mass.

Another early indicator that the contagious phenomenon has gone viral is if it is spotted in a fresh region of the network, not contiguous with the seed area. Because propagation via long-range ties is self-sustaining, the first occurrence can be a useful indicator that the contagious phenomenon has reached critical mass and will therefore continue to spread via every available tie, regardless of its range. The caveat is that a false positive can also be caused by the homophilous clustering of nodes with low thresholds (“early adopters”), but

even in the case where a jump simply indicates an expansion to a set of “early adopters” the consequent increase in the number of infected nodes in turn increases the likelihood of a true positive.

This analysis also has practical implications for marketing strategy. The greater the need for social reinforcement to persuade individuals to adopt an innovation, the larger the size of the initial region of local propagation required for the contagious phenomenon to go viral. Thus, in deciding where to launch an innovation, the proportion of highly susceptible nodes in the initial region is less important than the overall size, so long as the nodes are sufficiently susceptible that the contagious phenomenon can spread through short-range ties.

The need for a critical mass also carries implications for initial pricing. For simple contagion, it may be optimal to set prices initially high, in order to maximize profits from the most interested customers, early adopters. In contrast, for complex contagion, it is better to set prices initially low to improve the chances that the contagious phenomenon will reach critical mass.

We note these implications of the existence of a critical mass not as policy recommendations but as suggested directions for theoretical and empirical research. Our analysis assumes highly stylized topologies composed of nodes with homogenous attributes. Much more research is needed before we can have confidence in the predicted existence of a bifurcation point in the propagation of complex contagion. Residential neighborhoods, college dormitories, and soccer stadia may loosely resemble a regular lattice, and other empirical social networks have been shown to have degree distributions that approximate a power law. Yet other social networks have degree distributions that are more irregular than a lattice and less skewed than a power law. In addition, the nodes in empir-

ical networks have heterogeneous attributes, thresholds that vary between innovators, early adopters, and laggards [14, 37], and influence that varies between influentials, opinion leaders, and followers [59]. Moreover, these attributes may be homophilously clustered. These complications preclude the ability to use a set of formal results to confidently predict the critical mass in natural settings, and systematic empirical research is needed to see if the predicted existence of a bifurcation point is observed in empirical social networks. The important contribution of the present study is not that we have settled the question but quite the opposite – the predicted existence of a critical mass opens up an important direction for both theoretical and empirical research on why some contagion phenomena “go viral” and others die.

CHAPTER 4

CHAPTER 4: EMPIRICAL SOCIAL CONTAGION

In this chapter, I focus on a dataset of empirical contagion phenomena in an effort to validate the complex contagion model [21] holds up in a real-world setting. I use a dataset of photo tags from Flickr, where each tag is considered to be a potential contagion phenomenon spreading through the network of friendships between Flickr users. I first examine the distribution of thresholds over the set of tags. I then analyze the tags to determine which, if any, are complex contagion phenomena that have reached critical mass as described in Chapter 3, and find that at least some such tags exist. Finally, I show that simple natural language processing and statistical techniques can be used to narrow the search space from to a much smaller space of candidates that meet the criteria for complex contagion. These techniques make possible the efficient identification of complex contagion phenomena even in very large datasets of potentially contagious entities.

4.1 Introduction

In the previous chapters, I have described several models for the analysis of contagious phenomena diffusing on a social network. In this chapter, I turn to empirical contagion - actual behaviors, products and ideas that diffuse on real-world social networks. I use the analysis from previous chapters to guide my exploration of empirical contagion: the focus of this chapter will be to explore how abstract concepts like contagion threshold and critical mass manifest in a real-world setting.

The first step towards empirical analysis of contagious phenomena is their identification. In principle, a contagious phenomenon can be anything that spreads on a social network. This definition, however, is too broad, as it presents serious problems for identifying contagion. Some contagious phenomena change as they spread, and tracking their transformation and diffusion is a formidable problem in and of itself, which I leave to future work. Other phenomena do not change as they spread, but are difficult to trace, as there is no written record of their diffusion. Memes and behaviors that spread exclusively through verbal networks, for instance, are very difficult to track. Therefore, I focus on what I call empirical traceable social contagion: behaviors, products and ideas that spread throughout a social network, and leave a digital trace of their passage through the social ties between nodes in this network.

I can leverage this definition to apply theoretical concepts like threshold, perimeter growth, and critical mass, described in previous chapters, to empirical contagion phenomena.

It is important to note that these concepts do not translate perfectly to real-world applications: empirical contagion phenomena have no “hard threshold” beyond which every individual will adopt the behavior in question, and below which no individual will adopt it. Rather, the relationship between adoption likelihood and number of exposures to adopted individuals is probabilistic, with more exposures often (but not always) increasing the likelihood of adoption [11, 41].

Furthermore, even if we relax concepts like threshold to allow for stochastic adoption decisions, there remains a host of factors related to adoption that are not covered by the analysis in previous chapters. Threshold heterogeneity,

or its stochastic counterpart the heterogeneity of adoption likelihood functions, means that some individuals are much more likely to adopt the same contagious phenomenon given the same number of adopter friends as other individuals. Variations in interpersonal influence means that some individuals can induce their friends to adopt even with one exposure when the baseline likelihood of adoption remains low. Diffusion of some contagious phenomena via multiple media (including mass media) means that some nodes will be exposed to them even before any of their friends adopt the phenomenon in question. Some phenomena diffuse better across stronger ties, meaning that the strength of relationships in social media [33] has an effect of adoption likelihood between specific pairs of agents. Finally, people can decide to un-adopt a contagious phenomenon instead of sticking with it forever once adopted.

Despite these confounding factors, it remains possible to isolate a few and focus on their effect on adoption behavior. This chapter will follow recent research [11, 60] to focus on large-scale effects first, looking for patterns at the level of an entire network and/or across a large set of contagious phenomena. This focus “washes out” statistical noise and idiosyncratic effects of interpersonal relationships, preferences, etc. However, having examined adoption behavior at a large scale, we also investigate a few specific examples as case studies. Case studies produce results which may be more idiosyncratic, but give a much more detailed picture of the dynamics of, for instance, a particular contagious phenomenon. Future work can leverage the results from case studies to perform more large-scale analysis, completing the macro-micro analysis cycle.

4.2 Theory

Before moving on to analysis of empirical data, we first describe how the concepts of contagion threshold, first introduced by [34] and studied in detail by [49, 21], and contagion critical mass, introduced in the previous chapter, translate to an empirical setting. Consider an empirical traceable social contagion phenomenon c as described in the introduction to this Chapter. Then the digital trace of c 's passage through the social network consists of two datasets:

- A list $AT(c)$ of (node n , timestamp t) *adoption tuples* where each tuple indicates the adoption of c by n at t .
- A list $RT(c)$ of (node $n1$, node $n2$, timestamp t) *relationship tuples* where each tuple indicates the existence of an edge between $n1$ and $n2$ at t .

4.2.1 Threshold

In theory, the threshold of a contagious phenomenon is the critical number of infected neighbors k at or above which any node n will automatically adopt c (and below which no node n will ever adopt). In practice, such thresholds do not exist; however, for any given k , we can calculate the likelihood $p(k)$ that a node with k neighbors will adopt c . I follow [60] in the method of calculation: let $E(k)$ be the number of agents who have had at any point k friends who have adopted c , and $I(k)$ the number of agents who have both had k friends who have adopted c , and themselves adopted c before the $k + 1$ st of their friends did so. Then let $p(k) = I(k)/E(k)$. The latter work calls $p(k)$ calculated over a range of k

for a particular contagious phenomenon the exposure curve of c .¹

For contagion with a true threshold, the exposure curve is a step function, going from 0 below the threshold to 1 at and above the threshold. We can approximate the threshold of an empirical contagion phenomenon by identifying the region where the exposure curve most resembles a step function, and we can give a confidence score for our approximation by estimating the extent to which the exposure curve resembles a step function in that region. More formally:

Definition 4.2.1. The threshold $a(c)$ of an empirical contagion phenomenon c is:

$$a(c) = \operatorname{argmax}_{k \in [1 \dots K_{max}]} \Delta(p(k), p(k-1))$$

Definition 4.2.2. The confidence score $CS(a, c)$ for threshold k of an empirical contagion phenomenon c is:

$$CS(a, c) = \frac{\Delta(p(a), p(a-1))}{\text{AvgDelta}} - 1$$

where AvgDelta is the average $\Delta(p(k), p(k-1))$ for $k \in [1 \dots K_{max}] \wedge p(k) > p(k-1)$. The latter restriction is necessary so as not to deflate the average artificially over regions where $p(k)$ actually decreases. We can ignore those regions as not being candidates for points where $p(k)$ resembles a step function.

Intuitively, a step-like exposure curve will have a point where $p(k)$ grows much faster than in the rest of the region, so $a(c)$ will be the corresponding value

¹In practice, for high values of k , there may be very few nodes with that many friends, and $p(k)$ is very noisy, so we focus on some subset of values $k \in [0 \dots K_{max}]$

of k , and the confidence score will be high, as the growth rate may be several times the average. In contrast, a more gradual exposure curve may have a point where $p(k)$ grows a little faster than in the rest of the region (possibly due to noise), but the confidence score will be low, as the growth rate will be very close to the average.

4.2.2 Critical Mass

In the previous chapter, we discussed the notion of contagion perimeter and contagion growth rate. I repeat the definitions here:

Definition 4.2.3. Given a threshold a or z/n contagious phenomenon with I nodes already infected, the *perimeter* $x(I)$ of that infected set is the set of nodes that have a or more (or z/n or more) infected neighbors.

The growth rate of a contagious phenomenon over time is the a ratio of the number of nodes about to be infected (the perimeter) to the number of nodes already infected, expressed as a function of I .

Definition 4.2.4. The growth rate of a contagious phenomenon $\lambda(I)$ is given by $\frac{x(I)}{I}$ as a function of I .

For empirical contagion, I modify Definition 4.2.3 as follows:

Definition 4.2.5. Given a threshold a or z/n contagious phenomenon with I nodes infected at time t , the *perimeter* $x(t)$ of that infected set is the set of nodes become infected at time $t + 1$.

This new definition reflects both the fact that empirical contagion phenomena have no true threshold, and the fact that empirical contagion data contain explicit adoption timestamps in AT . We can now rewrite the Definition 4.2.4:

Definition 4.2.6. The growth rate of a contagious phenomenon $\lambda(t)$ is given by $\frac{x(t)}{Adopt(t)}$ as a function of t .

Where $Adopt(t)$ is the number of users who have adopted the contagious phenomenon at timestamp t . Also in the previous chapter, we described that the shape of the growth rate curve differs based on whether a contagious phenomenon has reached critical mass. For contagious phenomena with threshold greater than one, prior to critical mass, $\lambda(t)$ will be decreasing as $Adopt(t)^{-1}$, whereas post critical mass, $\lambda(t)$ will at first increase polynomially in $Adopt(t)$ until the contagious phenomenon exhausts the network, and then $\lambda(t)$ will decrease linearly in $Adopt(t)$. The achievement of critical mass corresponds to the only period of growth in $\lambda(t)$.

For an empirical contagion phenomenon, reaching critical mass may not lead to complete network takeover. the contagious phenomenon may reach a critical mass and spread rapidly, but then reach a region of the network where the threshold is much higher than elsewhere, and stop spreading. External factors such as the appearance of new contagious phenomena may divert the attention of the networked population from c and limit its rate of spread. At the same time, external events such as the broadcast of the contagious phenomenon through different media, may create a burst of growth for c even absent critical mass. Nevertheless, contagious phenomena whose growth rate does not increase are less likely to have reached critical mass than contagious phenomena whose growth rate increases. We therefore define the *criticality* $\rho(c)$ of contagion

c as follows:

Definition 4.2.7. The criticality $\rho(c)$ of a contagious phenomenon c is the average increase in the growth rate of c over the total timespan of adoption timestamps in $AT(c)$.

This definition encompasses both contagious phenomena that experience many short bursts of growth and one sustained long burst of growth. In further work, I hope to differentiate between these two dynamics.

A different approach to critical mass requires analyzing the average tie range of edges through which the contagious phenomenon spreads. As I discussed in Chapter 3, a critical mass of adopters for complex contagion is the point where the contagious phenomenon goes from only being able to spread via short-range ties to being able to spread via short- and long-range ties, as a simple contagion. Empirical complex contagion phenomena that reach critical mass should, therefore, exhibit a jump in the average range of ties through which they spread around the critical mass point. We define average tie range for empirical contagion phenomena as follows:

Definition 4.2.8. The average tie range $AT(c)$ of contagion c at time t is the average tie range between all users who have adopted the contagious phenomenon at time t and their neighbors who have adopted c at some time $t' < t$.

By mapping t to number of adopters at time t , it becomes possible to identify a critical mass point for some empirical complex contagion phenomenon c as the number of adopters where $AT(c)$ experiences a burst of growth.

4.3 Data and Methods

I now apply the above concepts of AT , RT , threshold and criticality to a specific set of empirical contagion phenomena: the spread of tags on the Flickr photo sharing service. Flickr is a social media tool that allows users to upload, share, and tag photos. Flickr users can befriend each other. The friendships are mutual relationships, so if a is a friend of b , then b is a friend of a . Flickr tags are single-word text metadata that can describe anything from the subject of the photo (“sunset”) to the camera used to take the photo (“nikon”) to the community or competition the photo was submitted to (“goldstaraward”).

I use a Flickr dataset that contains records of all photos uploaded between January 1st and July 1st, 2008, along with the upload timestamp and the tags applied to the photo.² The dataset contains 60 million million photos and over 1.6 million thousand tags. A separate Flickr dataset contains records of friendships between 500 thousand users in that time period. Unfortunately, Flickr friendships do not have associated timestamps, so I must assume that all friendships are present at all times during the dataset. This assumption introduces some noise into my analysis, however, social network structures generally change more slowly than the spread of contagious phenomena (for all but the most costly contagious phenomena, the cost of adding or removing a tie is greater than the cost of adoption, so an agent is in general more likely to adopt some contagious phenomenon than to change his or her social network), so it is possible that the noise won’t significantly distort measurements of threshold and critical mass across many contagious phenomena.

²It is important to note that photos can be tagged after upload, which introduces some noise into our adoption tuple records

I treat the tags in this dataset as potential contagious phenomena that can spread through the Flickr friendship network. For these tags, I first generate the associated datasets AT and RT . For each Flickr tag τ , $AT(\tau)$ contains tuples (u, t) where u is a user who tagged at least one of her photos with τ and t is the first timestamp associated with any such photo. At the same time $RT(\tau)$ contains tuples (u_1, u_2, t) that indicates that Flickr users, at least one of whom tagged at least one of her photos with τ were friends at some point between January 1st and July 1st, 2008.

Some of these tags may in truth be not contagious phenomena at all, but achieve popularity independently of the Flickr friendship network. For example, we can't assume that if Flickr user a uses the tag "sunset" and then her friend b uses the same tag, that a somehow infected b with the tag. Both users may just have seen a sunset that they found visually appealing. In contrast, tags like "nikon" represent the usage of a camera product that may have spread virally from one Flickr user to another.

Overall, the issue of confounds like external sources or user homophily that create contagion-like diffusion patterns of products that are not truly contagious, is complex and subtle. Recent work by Aral et al.[7] investigates the effect of homophily on creating such patterns. Not having access to user profile data at the granularity required to replicate Aral et al.'s methods, I adopt a simplistic solution at the level of all tags: remove from the set of all tags any that are common English words from a dictionary word frequency list. It is important to note that some tags that correspond to common English words may in fact be contagious (e.g. "apple" when it refers to the personal computer brand), it isn't possible to differentiate between uses of these words as common nouns

(picture of an apple) vs. as proper nouns (picture of an apple computer) with the Flickr dataset I use. When focusing on specific tags later in this chapter, I analyze the dynamics of tag diffusion further to look for signs of true contagious phenomena as opposed to contagion-like diffusion patterns.

It is also important to note that analysis of critical mass will help identify tags that are not true contagious phenomena. First, such phenomena should be characterized by an absence of rapid spurts in perimeter growth or average tie range, since they spread independently of the network. In contrast, they should spread in a smooth way, with perimeter and tie range changing at a near-constant rate throughout the diffusion period. I examine the perimeter growth and tie range curves of a few tags that are candidate contagion in detail in the Results section, to make sure this is not the case.

Finally, tags with very few adopters have very sparse AT data, so the resulting exposure curves and growth rate data may be too noisy to accurately determine contagion threshold and criticality as given above. Therefore, from the set of Flickr tags that are not common English words, I focus only on the top 500 tags by total number of adopters as of July 1st, 2008. This is the final set of tags T .

For each tag in $\tau' \in T$, I use $AT(\tau')$ and $RT(\tau')$ to generate exposure curves and growth rate curves, and I use those curves to estimate the tag's threshold and criticality, as described in the previous section.

4.3.1 Tag Coding

I also perform some manual coding of the tags τ . This coding is done for exploratory purposes and to replicate the recent work of [60], but was done entirely by myself. Multiple coders are necessary to validate the results. My coding reveals that the tags fall more or less neatly into one of four categories: location / subject (referring to the location or subject of the photo), technique (referring to the photographic technique or software used to create the photo), camera (in some ways this category is a subset of technique, but it refers specifically to the camera used to create the photo, and there are as many) and group (referring to the community or competition this photo belongs to).

There are 260 subject / location tags, including “baybridge” and “mediterraneo”; 44 technique tags including “depthoffield” and “backlight”; 21 camera tags including “fujifilm” and “canoneos400d”; and 106 group tags including “absolutelystunningscapes” and “catchycolors”. Of the latter two tags, the first refers to a community that collects photos of “stunning ‘scapes of all kinds - landscapes, waterscapes” and the second refers to a competition.

4.4 Results

4.4.1 Threshold

I begin by describing the distribution of thresholds and confidence scores for the 500 tags studied. Figure 4.1 shows the distribution of confidence scores $CS(a, c)$ by threshold value $a(c)$ for the top 500 tags. For sparsity reasons, I set $K_{max} = 6$.

Note that the threshold with the highest average confidence scores are 1 and K_{max} , respectively. The high confidence value for threshold 1 is due to the fact that $p(0)$ is typically very small. The denominator of $p(0)$ is the number of people who have 0 friends who have used the tag (which is the entire population), while the numerator of $p(0)$ is the number of people who began using the tag before any of their friends did (a subset of the entire set of tag users, which is usually a tiny percentage of the entire population). The very large denominator means that $p(0) \ll p(1)$ and that $[0, 1]$ is usually the region where $p(k)$ most resembles a step function. The high confidence values for threshold K_{max} is likely due to the increased noise in $p(k)$ for higher thresholds (note that the confidence interval in $CS(a, c)$ increases as (a, c) does. Setting K_{max} to higher values and rerunning the plot confirms that for values of $k > 6$, values $CS(a, c)$ are much noisier than for lower thresholds.

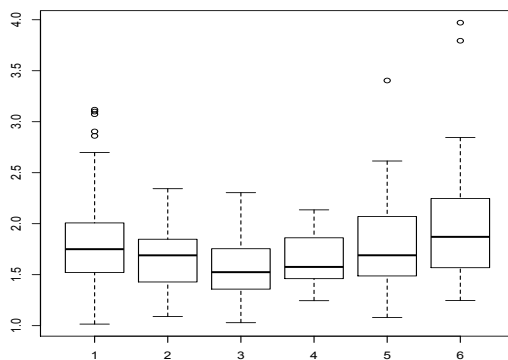


Figure 4.1: Confidence score distribution by threshold value for Flickr tags

Figure 4.2 shows the distribution of threshold values across manually labeled categories described in section 4.1. Note that the “camera” and “technique” categories have an average threshold of 1, while the “group” and “location” categories have a higher average threshold close to 2. This means that

technological tags on Flickr are much more likely to spread as simple threshold-1 contagion, whereas tags that are related to social aspects of Flickr like groups and the subject / location of the photo are much more likely to spread as higher-threshold contagion. This result is in line with the concept of complex contagion: while technological innovations like cameras and techniques can spread through weak ties and/or outside of Flickr (e.g. through advertisements), the social aspects of photography like groups and subject / location preferences are much more intrinsic to the social network of photographers, and more likely to spread via strong ties as friends engage in conversation on Flickr.³ The caveat is that, as seen in Figure 4.1, thresholds above 1 have lower confidence scores, so some tags in the groups and subject/location categories may in fact not be contagious phenomena at all, but become independently popular due to homophily or other effects.

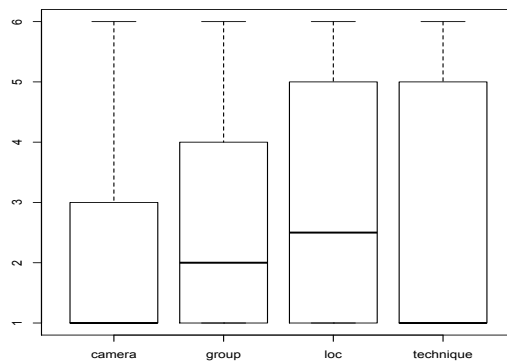


Figure 4.2: Distribution of threshold value by manually labeled category for Flickr tags

Figure 4.3 shows a histogram of the persistence parameter $F(P)$ over all 500 tags. Note that the distribution is close to normal, with a slight bump at the highest interval.

³conversation on Flickr is possible through comment streams on photos

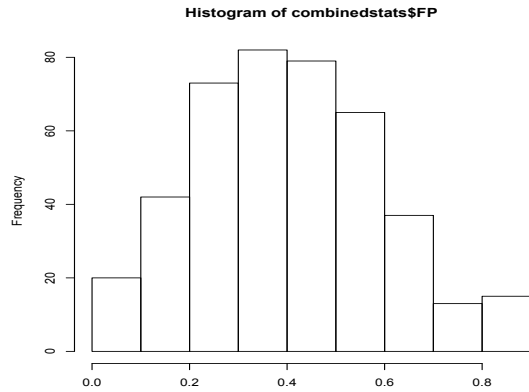


Figure 4.3: Histogram of persistence parameters for Flickr tags

4.4.2 Critical Mass

I now turn to an analysis of critical mass across the set T of Flickr tags. I begin with a plot of the criticality of each tag τ vs. the total number of adopters of τ in the dataset, shown in Figure 4.4. The two quantities are strongly linearly related (linear regression $R^2 = .81$, fitted line shown in red on plot). This is to be expected: the higher the criticality, the faster the contagious phenomenon grows, the more adopters it accumulates. Still, there are a few outliers, tags that have many more adopters than expected given their average perimeter growth rate. Looking at the data, these adopters turn out to be the tags “iflickr” (the Flickr app for the iPhone), “hdr” (the photography technique of high-dynamic range imaging), and “2008” (the tag corresponding to photos taken in 2008). Since the dataset covers the first half of 2008, the high growth rate of the last tag is not surprising and likely has nothing to do with contagion. We investigate the other two tags in the next section.

Figure 4.5 shows the distribution of criticality values by contagion threshold $a(c)$. Note that the criticality values (on the y axis) decreases slightly as threshold

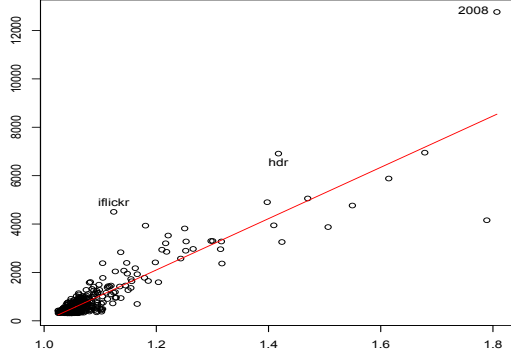


Figure 4.4: Contagion criticality $\rho(c)$ vs. number of adopters for Flickr tags

increases. This finding is in line with the theory of complex contagion, as it suggests that higher-threshold contagious phenomena grow more slowly. For $k = K_{max}$, the growth rate again rebounds, though that may be an artifact of noisy data.

On this graph as on the previous one, I identify several outliers - points with high criticality values and high thresholds. These turn out to be the tags “platinumphoto”, “naturesfinest”, “soe”, and “superbmasterpieces”, all referring to awards photo competitions on Flickr. Note that such competitions and awards can be contagious as a friend may nominate their friends’ photos to enter the competition or win the award. We investigate these tags in more detail in the next section.

Finally, I analyze critical mass by examining the dynamics of the average tie range $AT(t)$ vs. t for different tags. In order to compare multiple tags, I perform two transformations to the x axis of the plot: first, I map each value of t for a particular contagion c to the number of adopters $|Adopt_c(t)|$ who adopted c at t . Second, I divide $|Adopt_c(t)|$ by the total number of adopters $|Adopt_c|$ to transform

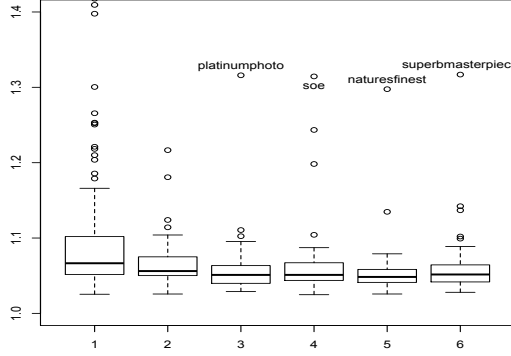


Figure 4.5: Contagion criticality $\rho(c)$ vs. threshold $a(c)$ for Flickr tags

a raw number of adopters into a percentile $Per_c(t)$ of the total adopter population. I then plot $AT(t)$ vs. $Per_c(t)$ for a number of different contagion. Due to the computational complexity of calculating tie range at a large scale, I constrain my analysis as follows: first, I only plot $AT(t)$ for the top 50 tags by overall number of adopters. Second, I transform tie range from a continuous integral variable to a categorical variable with values 2 (for nodes that are neighbors and have a neighbor in common), 3, and 4+. This limits very computationally expensive calculations of ties with ranges above 3, without overly skewing the results, as most ties will have a range of $O(\log(N))$, which is not much higher than 4 for the Flickr network. The final results are plotted in Figure 4.6.

This graph suggests two results: first, the general trend of tie range vs. percentile is smooth and positive. This trend indicates that most popular tags do spread over longer-range ties over time, but do not have a true critical mass point - there is no sudden jump in the tie range values. Second, there are a number of exceptions to this general trend. In Figure 4.6 I highlighted three contagion that go through a “bursty” period of tie range growth. These contagion are the “iflickr” tag (discussed in more detail below), the “smartphone”

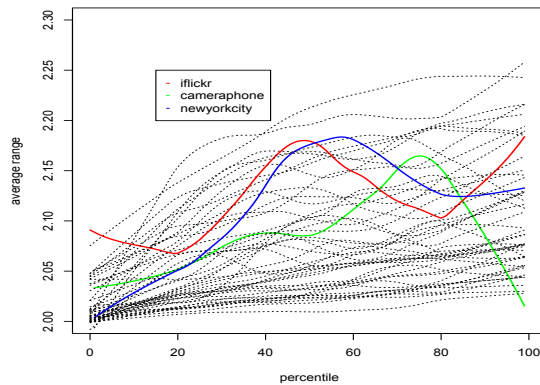


Figure 4.6: Average tie range for ties through which contagion spreads vs. percentile of adopters

tag (corresponding to pictures taken via smartphone) and the “newyorkcity” tag (corresponding to pictures taken in New York City). Each one of these tags goes through a period of rapid tie range growth, which suggests a possible critical mass point. For “flickr” and “newyorkcity”, the candidate critical mass point is around 20% of all adopters, whereas for “smartphone” the critical mass point is around 50% of all adopters.

These results are interesting in two ways. First, the candidate critical mass points for the three highlighted tags are larger (in terms of adopter percentage) than critical mass points for preferential attachment graphs theoretically derived in Chapter 3. Second, every one of these three tags goes through a burst period of growth followed by an *decrease* in the average tie range, suggesting that the tag goes through a period of spreading via short ties after it spreads via long ties. Clearly, network structure is not the only factor in determining critical mass. However, as more detailed analysis of the “flickr” tag in the next subsection suggests, we cannot attribute these fluctuations entirely to external events such as mass media coverage. It is most likely that the effect of network

structure is mediated by threshold heterogeneity, external stimuli, interpersonal influence and other factors with the final result showing dynamics qualitatively similar to, but quantitatively different from, the predictions of the complex contagion model.

The results by threshold and critical mass are in line with the theory of complex contagion and critical mass outlined in previous work. However, I do not find any clear clustering of threshold or criticality values - there is no easily identifiable set of tags that has a much higher threshold or average growth rate than others. Instead, I discover that the distributions of these quantities are skewed, with a high density of low-threshold, low-growth rate tags and a few outliers that score unusually high on one or both of these measures. Similarly, while I find that criticality has a strong linear relationship with total number of adopters, a few outliers have an unusually high number of adopters given their average growth rate. Thirdly, the tie range dynamics for a few tags are much more indicative of critical mass than the general trend. The next logical step is to investigate some of these outliers in more detail, to see what they tell us about the behavior of unusually fast-growing or popular tags.

4.5 Case Studies

I begin with a table of the six tags identified in the previous section.

Table 4.7 shows the threshold $a(c)$, the confidence $CS(a, c)$, the average growth rate $\rho(c)$, the total number of adopters, and the timespan (last timestamp - first timestamp of adoptions in the dataset, in days) for six outlier tags identified in the previous section. These tags represent every contagion threshold in

tag	a(c)	CS(a,c)	$\rho(c)$	totaladopters	timespan
iflickr	1	2.03	1.12	4504	170
hdr	1	1.76	1.42	6911	180
platinumphoto	3	1.43	1.32	3279	180
soe	4	1.78	1.31	2960	180
naturesfinest	5	1.23	1.30	3300	180
superbmasterpiece	6	1.34	1.32	2368	180

Figure 4.7: Threshold and criticality statistics for six contagion

the range except 2. Exploratory analysis of the six tags shows that “iflickr” has some of the most interesting dynamics, so below I focus on that tag. There is another reason to focus on “iflickr” - it is associated with the flickr application for the iphone, so by analyzing the spread of this tag, we can infer the spread of that application. This analysis has interesting implications for viral marketing, as it offers a useful proxy for application developers to study the spread of their product. At the end of the section, I present summary dynamics results for all six outlier tags for comparison to “iflickr.”

I begin by looking at the growth rate $\lambda(t)$ for iflickr in more detail. In Figure 4.8 I show the total number of adopters (orange dots) and $\lambda(t)$ (black line) as a function of time. We see that the growth rate experiences several jumps early on (probably due to noise, as the number of adopters is very small), and two jumps later on, one around day 75 and a second around day 110. These jumps both correspond to jumps in the number of adopters; furthermore, the slope of the number of adopters increases after the jumps.

Intuitively, these jumps indicate points where the iflickr tag may have reached critical mass and began to spread much faster. But what if they were caused by events external to flickr? To check for external influence on the diffu-

sion of the flickr tag, I first consulted the flickr blog to look for major releases or announcements around the days corresponding to the jumps (mid-March and end of April, 2008, respectively), and found none. I next used the Google Trends service to find points of major news / blog coverage for the string “flickr.” The black vertical lines in Figure 4.8 correspond to two such points - they both come about a week after the respective jumps in the growth rate curve. It seems that the growth spurts in tag adoption anticipate spurts in external coverage of the tag.

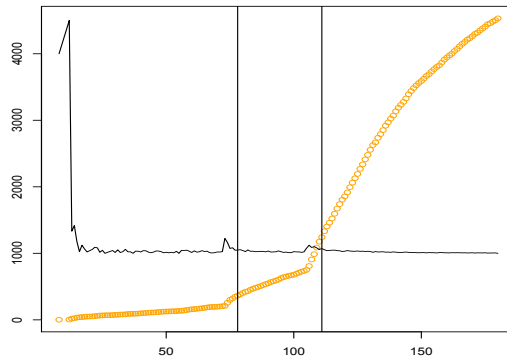


Figure 4.8: Number of adopters and contagion growth rate by time for the iFlickr tag

I next examine several threshold-related measures for flickr as a function of time. Figure 4.9 shows the total number of adopters (orange dots) as well as the number of Flickr users who have not yet adopted the flickr tag but have redundant (two or more) ties to tag adopters, as a function of time (green line). Vertical lines corresponding to Google trends data are repeated from Figure 4.8 for reference. Note that the number of redundant ties from adopters to non-adopters goes through several spurts, the first two of these corresponding to (and slightly anticipating) spurts in the number of adopters.

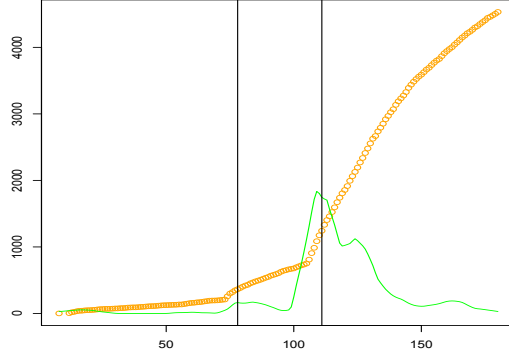


Figure 4.9: Number of adopters and number of redundant ties from adopters to non-adopters by time for iFlickr tag

Figure 4.9 tells us about number of nodes multiply exposed to iflickr adopters, but not about the actual adoption events. I next bin the adopters of iflickr by the timestamp during which they adopted, and for each bin calculate the average number of adopter friends $\bar{a}(a)$ each user in the bin had at that timestamp. This measure approximates the average threshold of adoption for Flickr users in the bin. The result is plotted in Figure 4.10: total number of adopters (orange dots), $\bar{a}(a)$ for each value of t (brown line), and $\bar{a}(a)$ for each value of t excluding adopters who had no adopter friends at the time (blue line).

The brown line essentially tracks the overall spread of iflickr through the Flickr population (regardless of network structure), whereas the blue line tracks only the adopters who may have decided to use their tags because of their friends' influence, and so the *potential* spread of iflickr through the Flickr friendship network. Figure 4.10 shows that at first both lines are relatively flat, but the brown line increases prior to the first jump in the number of adopters, and the blue line increases prior to the second such jump. The implications of these

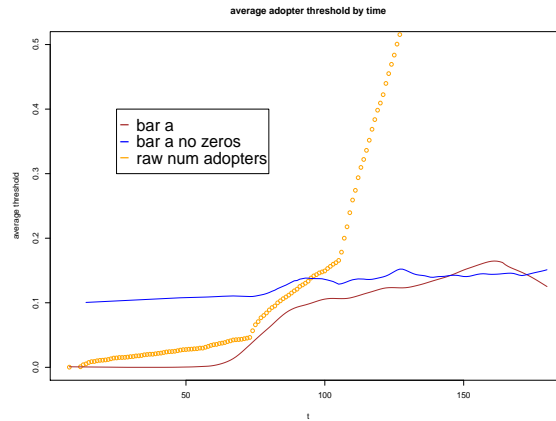


Figure 4.10: Average number of adopter friends by time for the iFlickr tag (brown - all tag adopters, blue - only tag adopters with one or more adopter friends)

two jumps are very interesting: at first, when both lines are relatively flat, and the brown line is near 0, the average threshold of iflickr adopters is also near 0. This implies that the tag is essentially spreading independently of the Flickr network. Shortly prior to the first jump, the brown line begins to grow, while the blue line stays relatively flat. In this region, more and more users with threshold greater than 0 begin to adopt iflickr; at the same time, users with threshold greater than 1 are not yet adopting iflickr (otherwise, the blue line would be growing). Finally, the blue line begins to grow prior to the second jump in the number of adopters. This region corresponds to the period where more and more users with threshold greater than 1 adopting flickr, as both lines are growing. In summary, Figure 4.10 suggests that the two jumps in the number of adopters are preceded by an increase in the average adopter threshold, which means the contagious phenomenon is expanding into populations that assign a higher cost to adoption.

Figures 4.9 and 4.10 are consistent with the notion of critical mass, as both

the number of exposures to multiple adopters and the number of adoptions given multiple exposures increase shortly prior to a rapid increase in the total number of contagion adopters. Therefore, plots of similar quantities for other tags could provide a useful litmus test for whether a contagious phenomenon is reaching critical mass (and, as a baseline, whether the contagious phenomenon is sprading through the network at all). For simplicity, we reproduce Figure 4.10 for the six outlier tags in Table 4.7. All six plots are on the same temporal scale ($t \in [0...180]$ days) on the x axis and the same abstract scale $[0...5]$ on the y axis. The results are shown in Figure 4.11.

Note that iflickr is the only tag where the average number of adopter friends is initially close to 0, and the only tag where the number of adopter friends for all adopters vs. adopters with at least one friend grow closer together. For “hdr”, these two lines are parallel though somewhat far apart. For the other four tags, the lines are parallel and very close together. Finally, note that the “superbmasterpiece” tag is the only one with a pronounced and long decline in average number of adopter friends, which precedes a marked decline in the contagious phenomenon growth rate.

4.6 Discussion and Conclusion

In this chapter, I’ve analyzed Flickr tags both at a large scale (a subset of 500 tags) and in detail (six specific tags). The results of the large scale analysis are consistent with Centola and Macy’s [21] theory of complex contagion, while the in-depth analysis suggests that a closer look at empirical contagion phenomena is needed to determine their dynamics with respect to critical mass.

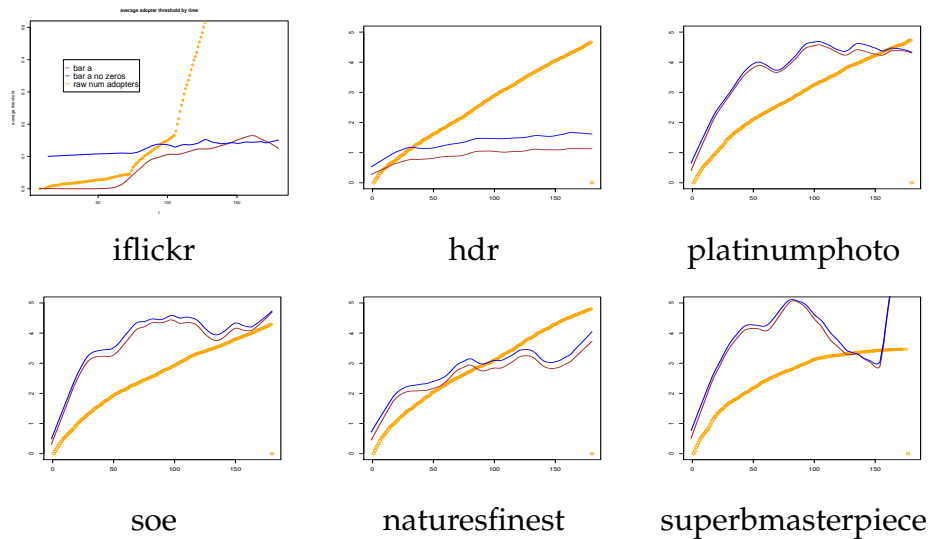


Figure 4.11: Average number of adopter friends by time for six different tags

While iFlickr has a very large number of adopters for its criticality $\rho(c)$, it also has a small threshold and a large confidence score in that threshold (relative to other outlier tags). Indeed, Figure 4.10 confirms that a substantive proportion of iFlickr adopters were the first of their friends to use the tag. However, the same figure also reveals that over time, the tag spreads to populations with a much higher adoption threshold.

This sort of fine-grained analysis allows one to make much more confident statements about the dynamics of individual contagion on Flickr, whether these contagion indeed spread through the social network, whether they achieve critical mass, and so on. At the same time, the large-scale analysis makes it possible to select a small number of interesting tags that could behave like contagious phenomena, without having to manually inspect the dynamics of hundreds or thousands of candidates. In future research, I hope to find an efficient way to summarize the adopter friends curves in Figures 4.10 and Figure 4.11 and

allow for efficient analysis of critical mass dynamics across many contagious phenomena. Nevertheless, the advantages of in-depth analysis (such as being able to investigate Google trends and other sources for external factors that may contribute to tag diffusion) make it an important step in determining what tags (and, in general, what products and behaviors) are truly examples of social contagion.

The fine-grained analysis also seems to suggest that the “iflickr” tags is the only clear example of a contagious phenomenon in the dataset that may have reached critical mass. This is not entirely surprising, as Centola and Macy’s[21] original research describes complex contagion as a set of rare and fragile phenomena. The vast majority of contagious phenomena either never reach critical mass, or have such a low threshold that the concept of critical mass is irrelevant to them (this applies to simple contagion). Furthermore, the search for complex contagion phenomena is complicated by the fact that a contagious phenomenon’s threshold may change over time and/or threshold heterogeneity in the underlying population (Figure 4.10 is in line with the latter effect). While more efforts must be made to identify complex contagion phenomena and critical mass regions, the low yield of the current study may be representative of future work. Searching for the rare high-threshold behavior that becomes popular is like looking for a needle in a haystack; conversely, searching for low-threshold behaviors or behaviors that never become popular is trivial, because they are the modal types of contagious behavior.

Finally, it is important to note that in several of the temporal plots, pre-critical-mass contagion properties (change in the average number of adopter friends) anticipate post-critical-mass contagion properties in and out of the net-

work (change in the accumulation of new adopters, appearance on Google trends). Such results appear in Figures 4.8 and 4.11. The implication of these results is that it may be possible to *predict* critical mass in empirical contagion phenomena by looking at certain pre-critical-mass contagion properties. Given the heterogeneity of contagious phenomena and the noise inherent in empirical data, it is far too early to make strong claims about contagion prediction. Nevertheless, this area of research demands further exploration and in the future it may be possible to integrate properties of the ties between adopters and non-adopters into a broader predictive model of contagion dynamics.

5.1 Summary of Results

This thesis explores social contagion from three different perspectives. In Chapter 2, I focused on the mechanism of local information and its effect on contagion dynamics. I showed that it is possible to formulate an optimization problem for agents trying to adopt a contagious phenomenon only when a critical fraction of the entire population has adopted, but having to rely on local information about the adoption states of their network neighbors. I formulate the optimization problem as minimizing the time an agent spends behaving suboptimally - being a contagious phenomenon adopter when a critical fraction of the entire population hasn't adopted, or not being a contagious phenomenon adopter when a critical fraction of the entire population has adopted.

Furthermore, I showed that for three common network models - the Poisson Random Graph, the Small World Graph, and the Preferential Attachment graph - it is possible to formulate strategies for agents as a combination of relying on their local network neighbors or querying the network at random, to approximately solve this optimization problem. For the Poisson Random Graph, most contagion either spread in logarithmic time or not at all, so reliance on local network neighbors produces a near-optimal strategy for nearly all combination of threshold values. For the Small World Graph, I showed that if the graph is an unrewired lattice (probability of edge rewiring $p = 0$), then querying the network at random k times produces a reduction of at most \sqrt{k} in the time an agent spends behaving suboptimally. Finally, for the preferential attachment graph, I

showed that for very low thresholds contagion spread in logarithmic time, but as threshold increases, there is a delay factor as the contagious phenomenon must first infect low-degree nodes before spreading to the hubs. Simulations show that this delay factor is linear in the threshold.

Chapter 3 focuses on a subset of social contagion called complex contagion. Complex contagion phenomena are interesting because they require multiple reinforcement to spread - an agent will only adopt a complex contagion phenomenon if $a > 1$ of her friends have already adopted. Previous work [21] has shown that the requirement of multiple reinforcement induces a drastically different dynamics for complex contagion than for simple contagion ($a = 1$). Complex contagion spreads on Small World Graphs in a fragile way: for a given Small World Graph and a given threshold a , a contagious phenomenon with that threshold may spread throughout the entire network starting with a small number of seed nodes, or stop spreading after reaching just a few nodes outside the seed set.

We demonstrate that there exists a bifurcation point in the spread of complex contagion – the critical mass. For simple contagion phenomena like information and disease, this bifurcation does not exist. Such phenomena can leverage long-range ties even with only one infected node, hence a single seed is sufficient to create a critical mass.

For complex contagion phenomena, in contrast, the growth process will have two phases separated by a very sharp transition. Initially, the contagious phenomenon can only spread locally, that is, via short-range ties. Once every node that is reachable via short-range ties is infected, propagation terminates if the level of infection remains sub-critical. However, if the region reachable via local

propagation is sufficiently large, the contagious phenomenon will reach critical mass and “go viral.” On this side of the bifurcation point, the contagious phenomenon can now spread via long-range ties that allow it to “jump” to a fresh area. This area of infection then rapidly expands through short-range ties.

The analysis also has an important theoretical implication for understanding why some contagion phenomena “go viral” and others do not. For simple contagion to escape the region of initial infection, it need only reach a node with a long-range tie or a hub that can broadcast the contagious phenomenon more widely. On a small world or undirected power law network, that is guaranteed to eventually occur, so long as the contagious phenomenon remains capable of passing from one node to another (e.g., there is no decline in infectiousness such as might happen in a news cycle). For a complex contagion to go viral, it must infect sufficient nodes that a long-range tie or susceptible hub can make a difference.

Finally, in Chapter 4, I explore contagion in the real world, focusing on a set of photo tags that propagate through the Flickr photo sharing network. I begin with large-scale analysis of the photo tags, and find two results in line with the concept of complex contagion: first, I find that tags that correspond to social aspects of Flickr, like groups, are much more likely to spread as higher-threshold contagion phenomena than technological tags that describe the camera or technique used to take the photograph. Centola and Macy [21] posit that many complex contagion phenomena, such as joining social movements or the spread of rumors, are social in nature. Second, I find that higher-threshold contagion phenomena grow more slowly over time, which is consistent with Centola and Macy’s [21] observation that complex contagion needs to spread through short-

range ties, which limits the contagious phenomenon's growth rate.

Next, I focus on a few tags in detail. Of special interest is the "iflickr" tag that corresponds to the spread of the Flickr app for the iPhone. Early on, the tag is picked up by a large number of threshold 0 adopters (people who use the tag before any of their friends do so), but its growth rate is slow. Only after higher-threshold adopters pick up the tag does the growth rate of "iflickr" take off. This result is anecdotal but points to the fact that empirical contagion phenomena can and do reach critical mass, and can't rely on early adopters to spread throughout the target population. Furthermore, the reaching of higher-threshold adopters for iFlickr precedes a spurt in the tag's growth rate. Similar analysis of five other tags shows that minor changes in the number of higher-threshold adopters do not precede changes in overall growth rate, but drastic increases precede growth spurts while drastic decreases in the number of higher-threshold adopters precede slowdowns. These results call for more investigation, but suggest that it may be possible to anticipate a critical mass point for a contagious phenomenon by analyzing the average threshold of adopters over time.

5.2 Broader Implications

The results of this thesis have broad implications in the areas for viral marketing as well as policymaking around social contagion like social movements and rumors. Chapter 2 shows that network structure and contagion threshold play an important role in optimizing adoption behavior given local information. This is a somewhat counterintuitive result, since understanding network structure and threshold seem to require global information about the contagious

phenomenon. However, it may be possible to estimate both from local information or from prior studies. For instance, policymakers wishing to prevent a social movement from spreading to a particular area may have a general understanding of the global structural properties of the network (degree distribution, connectivity) and may be able to estimate the threshold of the social movement cheaply in an experimental setting or a small-scale field study. They can then leverage this information to decide whether to isolate the target area from the rest of the network or to flood it with new connections, depending on whether the new connections increase the likelihood of infection or induce a higher level of neighborhood diversity to the target area.

Chapter 3 shows that the growth rate of contagion can change drastically over time. Prior to critical mass, the growth rate decays, but as soon as the contagious phenomenon is able to spread via long-range ties, the growth rate reverses and rapidly accelerates. This qualitative change in the rate of growth from negative to positive is a statistical signature of critical mass. The existence of a bifurcation point in the propagation of complex contagion has a potentially valuable practical implication for the ability to predict the eventual outcome at the early stages of a viral marketing campaign.

The analysis of complex contagion also has practical implications for marketing strategy. The greater the need for social reinforcement to persuade individuals to adopt an innovation, the larger the size of the initial region of local propagation required for the contagious phenomenon to go viral. Thus, in deciding where to launch an innovation, the proportion of highly susceptible nodes in the initial region is less important than the overall size, so long as the nodes are sufficiently susceptible that the contagious phenomenon can spread

through short-range ties.

The need for a critical mass also carries implications for initial pricing. For simple contagion, it may be optimal to set prices initially high, in order to maximize profits from the most interested customers, early adopters. In contrast, for complex contagion, it is better to set prices initially low to improve the chances that the contagious phenomenon will reach critical mass.

Finally, Chapter 4 applies the results from Chapter 3 in an empirical setting and presents important methodological considerations for the study of social contagion. The in-depth analysis suggests that a closer look at empirical contagion phenomena is needed to determine their dynamics with respect to critical mass. At the same time, the large-scale analysis makes it possible to select a small number of interesting tags that are potential instances of social contagion, without having to manually inspect the dynamics of hundreds or thousands of candidates. In the future, I hope to find more efficient ways of combining the macro-scale and the micro-scale approaches.

The fine-grained analysis in Chapter 4 yields few positive results, which carries important implications for future studies of complex contagion. The vast majority of contagious phenomena either never reach critical mass, or have such a low threshold that the concept of critical mass is irrelevant to them (this applies to simple contagion phenomena). Furthermore, the search for complex contagion phenomena is complicated by the fact that a contagious phenomenon's threshold may change over time and/or threshold heterogeneity in the underlying population. The search for examples of complex contagion and critical mass will require intensive data mining for a small yield. Nevertheless, the discovery of examples like the "iflickr" tag shows the importance of this

area of research: the prediction of critical mass in an empirical setting could be a great boon to marketers and policymakers alike.

5.3 Future Work

It is important to acknowledge the limitations of current research on social contagion. Chapters 2 and 3 are heavily theoretical, and the application of the analysis therein to empirical social contagion phenomena is not straightforward. This analysis assumes highly stylized topologies composed of nodes with homogenous attributes. Some empirical networks have degree distributions that don't approximate those of Poisson Random Graphs, Small World Graphs, or Preferential Attachment Graphs. In addition, the nodes in empirical networks have heterogenous attributes, thresholds that vary between innovators, early adopters, and laggards [14, 37], and influence that varies between influentials, opinion leaders, and followers [59]. Moreover, these attributes may be homophilously clustered. These complications preclude the ability to use a set of formal results to confidently predict the critical mass or optimize behavior under conditions of local information in natural settings. Much more empirical research is needed to make the transition from theoretical models to applied predictive systems in the area of social contagion.

In future work, I hope to address several natural extensions to the study of empirical social contagion phenomena: attribute and threshold heterogeneity, tie strength (which may increase or decrease the probability of contagion transmission via a particular network tie), and contagion confounds like homophily. The ultimate goal of this research is to create a comprehensive predictive model

for social contagion phenomena that would be able to incorporate theoretical results from this thesis along with more standard predictive mechanisms (e.g. machine learning) to forecast the rate of spread, equilibrium number of adopters, and other dynamic properties of socially infective behavior. I do not expect to be able to accomplish this goal on my own, and call for more research on empirical social contagion. Recent work such as [60] along this vein is especially promising and suggests that this area of study is fruitful with many opportunities for further research.

BIBLIOGRAPHY

- [1] L.A. Adamic and E. Adar. Friends and neighbors on the web. *Social Networks*, in press.
- [2] Eytan Adar, Li Zhang, Lada A. Adamic, and Rajan M. Lukose. Implicit structure and the dynamics of blogspace. In *Workshop on the Weblogging Ecosystem*, 2004.
- [3] R. Albert and A. L. Barabási. Statistical mechanics of complex networks. *Reviews of Modern Physics*, 74:47–97, 2002.
- [4] R. Albert, H. Jeong, and A. L. Barabási. Attack and error tolerance of complex networks. *Nature*, 406:378–382, 2000.
- [5] R.M. Anderson and R.M. May. *Infectious Diseases of Humans*. Oxford: Oxford University Press, 1991.
- [6] H. Andersson. Epidemic models and social networks. *Mathematical Scientist*, 24:128–147, 1999.
- [7] Sinan Aral, L. Muchink, and A. Sundararajan. Distinguishing influence based contagion from homophily driven diffusion in dynamic networks. In *Proc. Natl. Acad. Sci. USA*, volume 106, 2009.
- [8] R. Axelrod. An evolutionary approach to norms. *American Political Science Review*, 80:1095–1111, 1986.
- [9] R. Axelrod. The Dissemination of Culture: A Model with Local Convergence and Global Polarization. *J. Conflict Resolut.*, 41(2):203–226, 1997.
- [10] R. Axelrod, W. Mitchell, R.E. Thomas, D.S.Bennett, and E.Bruderer. Coalition formation in standard-setting alliances. *Management Science*, 41:1493–1508, 1995.
- [11] Lars Backstrom, Dan Huttenlocher, Jon Kleinberg, and Xiangyang Lan. Group formation in large social networks: Membership, growth, and evolution. In *In Proc. 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2006.
- [12] A. L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.

- [13] A. Barrat and M. Weigt. On the properties of small-world networks. *European Physics Journal B*, 13:547–560, 2000.
- [14] John Berry and Ed Keller. *The Influentials: One American in Ten Tells the Other Nine How to Vote, Where to Eat, and What to Buy*. NY: Simon and Schuster, 2003.
- [15] H. Bott. Observation of play activities in a nursery school. *Genetics Psychology Monographs*, 4:44–88, 1928.
- [16] A. Broder. Graph structure in the web. *Computer Networks*, 33(1-6):309–320, June 2000.
- [17] R.S. Burt. General social survey network items. *Connections*, 8:119–123, 1985.
- [18] Campbell, Converse, Miller, and Stokes. *The American voter*. New York: John Wiley and Sons, Inc., 1960.
- [19] K. Carley. A theory of group stability. *American Sociological Review*, 56:331–354, 1991.
- [20] Damon Centola. The spread of behavior in an online social network experiment. *Science*, 329:1194–1197, 2010.
- [21] Damon Centola and Michael Macy. Complex contagions and the weakness of long ties. *American Journal of Sociology*, 113:702–734, 2007.
- [22] Michael Suk-Young Chwe. Structure and strategy in collective action. *American Journal of Sociology*, 105:128–155, 1999.
- [23] R. Cohen, D. ben Avraham, and S. Havlin. Efficient immunization of populations and computers. *Preprint cond-mat*, 2002.
- [24] J.S. Coleman. *Community Conflict*. New York: Free Press., 1957.
- [25] David J. Crandall, Dan Cosley, Daniel P. Huttenlocher, Jon M. Kleinberg, and Siddharth Suri. Feedback effects between similarity and social influence in online communities. In *14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2008.

- [26] de Solla. Networks of Scientific Papers. *Science*, 149(3683):510–515, July 1965.
- [27] K.W. Deutsch. *Nationalism and social communication: An inquiry in the foundations of nationality*. Cambridge, MA: Technology Press, 1953.
- [28] K.W. Deutsch. *Nationalism and its alternatives*. New York: Knopf, 1969.
- [29] O.D. Duncan, D.L. Featherman, and B. Duncan. *Sociometric Background and Achievement*. New York: Seminar, 1972.
- [30] Nathan Eagle, Michael Macy, and Robert Claxton. Network diversity and economic development. *Science*, 328:1029–1031, 2010.
- [31] P. Erdős and A. Rényi. On the evolution of random graphs. *Publications of the Mathematical Institute of the Hungarian Academy of Sciences*, 5:17–60, 1960.
- [32] Galeotti. Network games. *Review of Economic Studies*, forthcoming.
- [33] Eric Gilbert and Karrie Karahalios. Predicting tie strength with social media. In *27th International Conference on Human Factors in Computing Systems*, 2009.
- [34] Mark Granovetter. Threshold models of collective behavior. *The American Journal of Sociology*, 83(6):1420–1443, 1978.
- [35] P. Grassberger. On the critical behavior of the general epidemic process and dynamical percolation. *Mathematical Biosciences*, 63:157–172, 1983.
- [36] H.W. Hethcote. The mathematics of infectious diseases. *SIAM Review*, 42:599–653, 2000.
- [37] Elihu Katz and Paul Felix Lazarsfeld. *Personal Influence: the Part Played by People in the Flow of Mass Communications*. Transaction Publishers, 1955.
- [38] Michael L. Katz and Carl Shapiro. Network externalities, competition, and compatibility. *American Economic Review*, 75:424–440, 1985.
- [39] Ravi Kumar, Jasmine Novak, Prabhakar Raghavan, and Andrew Tomkins. Structure and evolution of blogspace. *CACM*, 47:35–39, 2004.

- [40] P. Lazarsfeld and R.K.Merton. Friendship as a social process: A substantive and methodological analysis. In *Freedom and Control in Modern Society*, pages 18–66. New York: Van Nostrand, 1954.
- [41] Jure Leskovec, Lada Adamic, and Bernardo Huberman. The dynamics of viral marketing. In *Proc. 7th ACM Conference on Electronic Commerce*, 2006.
- [42] Jure Leskovec, Mary Mcglohon, Christos Faloutsos, Natalie Glance, and Matthew Hurst. Cascading behavior in large blog graphs. In *In SDM*, 2007.
- [43] D.K. Lewis. *Convention: A philosophical study*. Cambridge, MA: Harvard University Press, 1967.
- [44] Seymour Martin Lipset, Martin Trow, and James S. Coleman. *Union Democracy: The Internal Politics of the International Typographical Union*. New York: Free Press, 1956.
- [45] C.P. Loomis. Political and occupational cleavages in a hanoverian village. *Sociometry*, 9:316–333, 1946.
- [46] A. Madanipour, G. Cars, and J. Allen (eds.). *Social Exclusion in European Cities: Processes, Experiences, and Responses*. London: Jessica Kingsley, 1998.
- [47] Gerald Marwell and Pamela Oliver. *The Critical Mass in Collective Action*. Cambridge University Press, 1993.
- [48] Stanley Milgram. The small world problem. *Psychology Today*, 2:60–67, 1967.
- [49] Stephen Morris. Contagion. *Review of Economic Studies*, 67:57–78, 2000.
- [50] Stephen Morris and Hyun Song Shin. Heterogeneity and uniqueness in interaction games. In *Cowles Foundation Discussion Papers*. Cowles Foundation for Research in Economics, Yale University, New Haven, CT, 2003.
- [51] R.R. Nelson and S.G. Winter. *An evolutionary theory of economic change*. Cambridge, MA: Harvard University Press, 1982.
- [52] M. E. J. Newman. The structure and function of complex networks. *SIAM Review*, 45(2):167–256, 2003.

- [53] M.E.J. Newman, S.H. Strogatz, and D.J.Watts. Random graphs with arbitrary degree distributions and their applications. *Physica Review E*, 64, 2001.
- [54] R. Pastor-Satorras and A. Vespignani. Epidemic dynamics and endemic states in complex networks. *Physica Review E*, 63, 2001.
- [55] R. Pastor-Satorras and A. Vespignani. Epidemic spreading in scale-free networks. *Physical Review Letters*, 86:3200–3203, 2001.
- [56] R. Pastor-Satorras and A. Vespignani. Immunization of complex networks. *Physica Review E*, 65, 2002.
- [57] R.D. Putnam. Political attitudes and the local community. *American Political Science Review*, 60:640–654, 1966.
- [58] Anatol Rapoport. Contribution to the theory of random and biased nets. *Bulletin of Mathematical Biology*, 19(4):257–277, Dec 1957.
- [59] Everett M. Rogers. *Diffusion of Innovations*. New York: Free Press, 5 edition, 2003.
- [60] Daniel M. Romero, Brendan Meeder, and Jon Kleinberg. Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on twitter. In *WWW 2011*, in press.
- [61] G. Saloner and J. Farrell. Installed base and compatibility: Innovation, product preannouncements and predation. *American Economic Review*, 76:940–955, 1986.
- [62] T. Schelling. Dynamic models of segregation. *Journal of Mathematical Sociology*, 1(1):143–186, 1971a.
- [63] J. Scott. *Social Network Analysis: A Handbook*. London: Sage Publications, 2 edition, 2000.
- [64] H. Simon. On a class of skew distribution functions. *Biometrika*, 42(3-4):425–440, 1955.
- [65] S.N.Dorogovtsev, A.V.Goltsev, and J.F.F. Mendes. Critical phenomena in complex networks. *Review of Modern Physics*, 80:1275, 2008.

- [66] John C Turner, Michael Hogg, Penelope J Oakes, and Stephen D Reicher. *Rediscovering the social group: a self-categorization theory*. Oxford: Basil Blackwell, 1988.
- [67] E. Ullmann-Margalit. *The emergence of norms*. Oxford, UK: Clarendon, 1977.
- [68] T.W. Valente. *Network Models and Methods for Studying the Diffusion of Innovations*. Cambridge University Press, 2005.
- [69] L.M. Verbrugge. The structure of adult friendship choices. *Social Forces*, 56:576–597, 1977.
- [70] A. Walker and C. Walker. *Britain Divided: The Growth of Social Exclusion in the 1980's and 1990's*. London: CPAG Ltd., 1997.
- [71] Duncan J. Watts. A simple model of global cascades on random networks. In *Proc. Natl. Acad. Sci. USA*, volume 99, pages 5766–5771, 2002.
- [72] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393:440–442, 1998.