

## **HKUST SPD - INSTITUTIONAL REPOSITORY**

Title	Efficient Estimation of Network Games of Incomplete Information: Application to Large Online Social Networks
Authors	Chen, Xi; Van der lans, Ralf; Trusov, Michael
Source	Management Science, v. 67, (12), December 2021, p. 7575-7598
Version	Published Version
DOI	10.1287/mnsc.2020.3885
Publisher	INFORMS
Copyright	Management Science. Copyright © 2021 The Author(s). https://doi.org/10.1287/mnsc .2020.3885, used under a Creative Commons Attribution License CC BY 4.0
License	CC BY 4.0

This version is available at HKUST SPD - Institutional Repository (https://repository.ust.hk/ir)

If it is the author's pre-published version, changes introduced as a result of publishing processes such as copy-editing and formatting may not be reflected in this document. For a definitive version of this work, please refer to the published version.



### Efficient Estimation of Network Games of Incomplete Information: Application to Large Online Social Networks

#### Xi Chen,<sup>a</sup> Ralf van der Lans,<sup>b</sup> Michael Trusov<sup>c</sup>

<sup>a</sup> Department of Marketing Management, Rotterdam School of Management, Erasmus University, 3062 PA Rotterdam, Netherlands; <sup>b</sup> Department of Marketing, Hong Kong University of Science and Technology, Clear Water Bay, Kowloon, Hong Kong; <sup>c</sup> Department of Marketing Robert H. Smith School of Business, University of Maryland, College Park, Maryland 20742 **Contact:** chen@rsm.nl, **(**) https://orcid.org/0000-0002-6260-1990 (XC); rlans@ust.hk, **(**) https://orcid.org/0000-0002-7726-8238 (RvdL); mtrusov@umd.edu (MT)

Received: May 29, 2017 Revised: February 5, 2019; October 9, 2019; June 9, 2020; September 18, 2020 Accepted: September 29, 2020 Published Online in Articles in Advance: June 30, 2021

https://doi.org/10.1287/mnsc.2020.3885

Copyright: © 2021 The Author(s)

Abstract. This paper presents a structural discrete choice model with social influence for large-scale social networks. The model is based on an incomplete information game and permits individual-specific parameters of consumers. It is challenging to apply this type of models to real-life scenarios for two reasons: (1) The computation of the Bayesian–Nash equilibrium is highly demanding; and (2) the identification of social influence requires the use of excluded variables that are oftentimes unavailable. To address these challenges, we derive the unique equilibrium conditions of the game, which allow us to employ a stochastic Bayesian estimation procedure that is scalable to large social networks. To facilitate the identification, we utilize community-detection algorithms to divide the network into different groups that, in turn, can be used to construct excluded variables. We validate the proposed structural model with the login decisions of more than 25,000 users of an online social game. Importantly, this data set also contains promotions that were exogenously determined and targeted to only a subgroup of consumers. This information allows us to perform exogeneity tests to validate our identification strategy using community-detection algorithms. Finally, we demonstrate the managerial usefulness of the proposed methodology for improving the strategies of targeting influential consumers in large social networks.

History: Accepted by Matthew Shum, marketing.

**Open Access Statement:** This work is licensed under a Creative Commons Attribution 4.0 International License. You are free to copy, distribute, transmit and adapt this work, but you must attribute this work as *"Management Science*. Copyright © 2021 The Author(s). https://doi.org/10.1287/mnsc.2020.3885, used under a Creative Commons Attribution License: https://creativecommons.org/licenses/by/4.0/."

Supplemental Material: The online appendix is available at https://doi.org/10.1287/mnsc.2020.3885.

Keywords: social influence • network games • Bayesian estimation of games • community detection • targeting • online social networks

#### 1. Introduction

With the development of information technology and, especially, the proliferation of different forms of social media sites, social influence has become an essential factor of consumers' decision making. The existence of social influence implies that choices are mutually dependent, such that the choice of a consumer both influences and depends on choices of other consumers in the social network (Hogan et al. 2003, Gupta et al. 2006, Gupta and Mela 2008, Ahn et al. 2015). The ability to identify and quantify the strength of social influence on choice within a network, thus, has numerous applications. For example, in the field of customer-relationship management (e.g., Kumar and Reinartz 2012), managers may want to focus greater retention efforts on consumers whose decisions to leave the firm would have stronger spillover effects on other customers (Ascarza et al. 2017).

In targeted marketing, choices by influential consumers might be considered as more profitable for marketing communications and promotional activities under the presumption that their choices would trigger desired actions by other consumers (e.g., product adoption or higher engagement in social networks; see Hinz et al. 2011, Toker-Yildiz et al. 2017, and van der Lans et al. 2010). In education, students may be more likely to study if their friends also study hard and have less time for other activities (Glaeser and Scheinkman 2001). Similarly, decisions about participating in crime, smoking, and the decision to join the labor force are all marked by social influence (Glaeser and Scheinkman 2001, Lee et al. 2014, Nicoletti et al. 2018).

Although the benefits of knowing the extent of social influence in the choices of consumers are intuitively appealing, the assessment task is challenging. 7576

First, identifying social influence from choice data is difficult because of confounding factors, such as correlated (homophily) and exogenous effects (Manski 1993). For example, the observation that two consumers make similar choices is not necessarily attributable to social influence. The behavioral similarities can be attributed to common traits (correlated effects or homophily) or the exposure to the same exogenous shocks (exogenous effects). Second, although game-theoretic models have become an essential tool for the analysis of social interactions (Brock and Durlauf 2001), the empirical application of these models to large social networks is limited. This is because the application faces a scaling problem (Hartmann 2010) that makes it difficult to apply such models to large social networks. For example, the extant algorithms such as nested fixed points and Mathematical Programming with Equilibrium Constraints have been used to study interactions of up to a few hundred people or firms (e.g., Zhu and Singh 2009, Hartmann 2010, Misra 2013, and Lee et al. 2014). However, these methods do not scale well to business settings that may involve interactions between tens of thousands of heterogeneous consumers.

To fill this gap in the literature, this paper extends the discrete choice model with social influence introduced by Brock and Durlauf (2001) and extended in Lee et al. (2014) to empirical settings with a large social network and individual-specific parameters. Allowing for individual-specific parameters is important, as consumers may respond differently to social influence, leading to different policy implications, as highlighted in our empirical applications. To address the scalability and identification issues as discussed above, we propose two novel solutions. First, we derive the uniqueness and local stability conditions of the game-theoretic model, which enable us to apply the scalable stochastic Bayesian estimation procedure introduced by Imai et al. (2009). By doing so, we expand the application domain of this pseudo-fixed-point procedure from structural demand models (Imai et al. 2009, Sun and Ishihara 2018) to games on large social networks. Second, to facilitate the identification of social influence, we propose to use community-detection algorithms (Fortunato 2010) to uncover latent communities that are likely to be exposed to different time-varying unobserved shocks. With the learnt communities, we construct excluded variables using community-specific time fixed effects. One of the key benefits of this novel identification strategy is that it does not require observing excluded variables that are often difficult to obtain in practice (Imbens 2014).

We apply the proposed methods to an online gaming data set that contains a complete network with login activities of over 25,000 users across 30 days. Common for many other online games, tasks in the game require team efforts, making social influence an important driving factor for users' login decisions. Importantly, this data set contains observed exogenous shocks (i.e., targeted promotions on holidays), which allow us to validate the identification strategy based on latent communities. In particular, we follow Rivers and Vuong (1988) and perform exogeneity tests for the structural models. The results of the tests support the value of communities for the identification of social influence. Based on the estimates of the structural model, we further quantify the social influence of a user by comparing the observed scenario to a counterfactual scenario where a user becomes inactive. We then apply this insight to policy simulations to illustrate the effectiveness of different targeting strategies. These policy simulations show the benefits of incorporating individual-specific parameters into the model, as it increases the expected performance of targeting strategies about 10-fold. Moreover, as a novel insight, our counterfactual studies illustrate that optimal social targeting strategies should not only consider whom to target, but also when to target certain consumers.

The remainder of this paper is organized as follows. In Section 2, we describe our proposed model, derive the local stability condition of the network game, and characterize its uniqueness. In Section 3, we discuss model identification and community-detection algorithms, introduce our estimation method, and demonstrate its scalability with simulation studies. In Section 4, we describe the data of our empirical application, which consists of users' login decisions in an online game. This section also reports estimation results and validation of our identification strategy, and provides policy simulations that demonstrate the managerial usefulness of the proposed model for targeting decisions. In Section 5, we conclude and discuss directions for future research.

#### 2. A Local Interaction Choice Model on Social Networks 2.1. The Model

Consider a group of *N* consumers connected by a social network. An  $N \times N$  adjacency matrix *A* is used to represent this social network, with nondiagonal elements  $a_{ij} \in [0, \infty)$  indicating the weight consumer *i* puts on consumer *j* when making a decision. The diagonal elements  $a_{ii}$ , by convention, are set to zero, and we allow connections to be asymmetric  $(a_{ij} \neq a_{ji})$ . Suppose these i = 1,...N consumers are making a binary choice decision at time period *t* by choosing one of the actions from a discrete set with  $d_{it} = \{0,1\}$ . Without loss of generality, the utility of choosing action  $d_{it} = 0$  is normalized to 0.

To specify the utility function of the discrete choice model with social influence, we follow Brock and Durlauf (2001) and Lee et al. (2014), but allow parameters to be individual-specific. Consequently, the latent utility of consumer *i* choosing  $d_{it} = 1$  at time *t* is assumed to be additive-separable and partitioned into three terms as follows:

$$U(d_{it} = 1 | x_{it}, d_{-it}, \Theta_i) = V(x_{it}; \beta_i) + S(d_{it}, d_{-it}; \gamma_i) + \varepsilon_{it}.$$
(1)

In this equation,  $x_{it}$  is a  $K \times 1$  vector of individual characteristics that may change over time, and  $\beta_i$ the associated parameter vector and assumed to be individual-specific. Parameter  $\gamma_i$  captures the susceptibility to social influence of consumer *i*. The term  $V(x_{it}; \beta_i)$  represents the deterministic private utility that depends on  $x_{it}$ .  $S(d_{it}, d_{-it}; \gamma_i)$  captures the social utility which depends on decisions of other consumers  $d_{-it}$ . The parameter set  $\Theta_{i}$ , which contains  $\{\beta_i, \gamma_i\}$  and possible parameters of the distribution of  $\varepsilon_{it}$ , is assumed to be compact. Finally,  $\varepsilon_{it}$  is an independent and identically distributed (i.i.d.) (across consumers and time periods) random shock and assumed to follow a cumulative distribution function *F* that is twice continuously differentiable and a probability distribution function f that is bounded and strictly positive.

To further specify the model, we assume that only the direct friends  $j \in N_i$  of consumer *i* will directly influence the decision of consumer *i*, with  $N_i$  representing all consumers that are connected to consumer *i* (i.e.,  $a_{ij} > 0$ ). Furthermore, we specify the social utility function as responding to the decisions of friends:<sup>1</sup>

$$S(d_{it} = 1, d_{jt}; \gamma_i) = \gamma_i \sum_{j \in N_i} a_{ij} \mu^e_{ijt} (d_{jt} = 1).$$
(2)

In Equation (2),  $\mu_{ijt}^e(d_{jt} = 1) = E_i(d_{jt} = 1|I_{it})$  is the expectation of the decisions of consumer *j* formed by consumer *i*. The expectation is formed based on the information set  $I_{it}$ , which includes the individual characteristics  $\{x_{jt}, \forall j = 1, \dots, N\}$  and social network structure *A*. The information set  $I_{it}$  is assumed to be publicly known.

For the data-generating process, we make the following assumptions: (1) The observed social network is mature, which implies that consumers do not form new connections or dissolve old ones, and there are no new consumers joining the network;<sup>2</sup> (2) consumers are myopic, which implies that observations of different time periods are independent realizations of the network game; and (3) the distribution of random shocks *F* is common knowledge, but the random shock  $\varepsilon_{it}$  is only observed by consumer *i*, and not by other consumers. The last assumption leads to a static game with incomplete information, where consumers react to the expectations over behaviors of their friends. Incomplete information is realistic in large networks, as consumers are likely to rely on incomplete or partial information (e.g., characteristics of their friends) to form expectations over friends' decisions (Lee et al. 2014). In contrast, previous research (e.g., Hartmann 2010) with the complete-information assumption considered social influence in small groups, where consumers can exchange information freely and frequently. Furthermore, a disadvantage of the completeinformation assumption is that it almost always leads to multiple equilibria (Galeotti et al. 2010), which limit its application to relatively small groups (Soetevent and Kooreman 2007). Finally, the literature on global game theory illustrates how the introduction of (a small amount of) incomplete information may lead to unique equilibrium solutions and provides interesting and novel economic intuitions (Carlsson and Damme 1993).

Given the specification above, the probability of consumer *i* choosing an action  $d_{it}$  at time *t* can be represented as,

$$\Pr(d_{it}=1|I_{it},\mu_{it}^{e}) = F\left(V_{it}+\gamma_{i}\sum_{j\in N_{i}}a_{ij}\mu_{ijt}^{e}(d_{jt}=1)\right).$$
 (3)

In Equation (3),  $V_{it} = V(x_{it};\beta_i)$ , the intrinsic utility term as specified above, and  $\mu_{ijt}^e$  is consumer *i*'s expectation about the behaviors of a friend *j*. To make the model tractable, we further impose the rational expectation condition (Hansen et al. 1991),<sup>3</sup> which implies that, in equilibrium, the expectations of consumers are consistent with the actual choice behaviors, as described by the equilibrium choice probabilities—that is,  $u_{ijt}^e d_{jt} = 1$ ) =  $p_{jt}^* (d_{jt} = 1)$ ,  $\forall i \in N \setminus j$ . Consequently, in equilibrium, the probability of consumer *i* choosing  $d_{it} = 1$  at time *t* is,

$$p_{it}^* = F\left(V_{it} + \gamma_i \sum_{j \in N_i} a_{ij} p_{jt}^*\right).$$
(4)

By stacking the equations above across all consumers, we obtain the following system of *N* equations:

$$\boldsymbol{P}_t^* = F(\boldsymbol{\Delta}_t + \boldsymbol{\Gamma} \boldsymbol{A} \boldsymbol{P}_t^*), \tag{5}$$

where  $\Delta_t = [V_{1t}, ..., V_{Nt}]'$  is an  $(N \times 1)$  vector containing consumers intrinsic utilities,  $\Gamma$  is an  $(N \times N)$  diagonal matrix, with  $\gamma_i$  as the *i*th diagonal element, and  $P_t^*$  an  $(N \times 1)$  vector containing equilibrium choice probabilities at time *t*. Combining Equations (1)–(5) and stacking the decision variables  $d_{it}$  in the  $(N \times T)$ matrix *Y* and  $x_{it}$  in the  $(N \times K \times T)$ -matrix *X*, we obtain the likelihood:

$$L(\Theta|Y, X, A) = \prod_{i \in Nt \in T} (p_{it}^*)^{d_{it}} (1 - p_{it}^*)^{1 - d_{it}},$$
  
subject to :  $P_t^* = F(\Delta_t + \Gamma A P_t^*), \forall t \in T.$  (6)

#### 2.2. The Problem of Multiple Equilibria

The applicability of the proposed game-theoretic choice model depends on whether there exists an equilibrium choice outcome given the individual-specific parameters. Under the assumption that *F* is continuous, there exists a solution to the system of Equations (5), based on Brouwer's fixed-point theorem. Although the system of Equations (5) admits an equilibrium solution, this solution may not be unique, and multiple equilibria could exist, as discussed in previous research (Brock and Durlauf 2001, Mazzeo 2002, Bajari et al. 2010, Ellickson and Misra 2011, Vitorino 2012). The existence of multiple equilibria would prevent us from pinning down a unique outcome and, thus, reduces the usefulness of our model for counterfactual studies that require equilibrium computations. For homogeneous parameters and a general network structure, there is no simple way to identify the number of equilibria (Lee et al. 2014). In the situation of individual-specific parameters, it becomes even more challenging to determine the number of equilibria. Fortunately, it is possible to derive the sufficient conditions for a locally stable equilibrium, which we further prove to be unique. This is particularly beneficial in the current situation of large social networks, which requires an efficient estimation method. The local stability condition theoretically justifies the focus on the unique equilibrium, which is necessary to apply the proposed estimation procedure. Next, we discuss these conditions.

#### 2.3. Ex Ante Equilibrium Selection with Local Stability

A commonly used method to deal with multiple equilibria is the "refinement of Nash equilibrium" (Myerson 1978). Following Brock and Durlauf (2001), we utilize a refinement criterion named "local stability." Local stability implies that the game converges to the original equilibrium if there is an infinitesimal perturbation in one or more consumers' equilibrium behaviors. Using such a refinement scheme has important theoretical and practical implications. First, if an equilibrium is not stable, a slight deviation by one of the consumers will cause other consumers to deviate further away from the equilibrium. In such situations, it is, thus, practically difficult, if not impossible, for the system to reach the equilibrium. Second, an important advantage of structural models is its ability to perform counterfactual analysis. Without the equilibrium refinement, we need to compute all possible equilibria, which might lead to multiple and possibly contradicting policy implications. Finally, the search for multiple equilibria becomes a nearly impossible task, with many heterogeneous consumers with individualspecific parameters.

Chen, van der Lans, and Trusov: Efficient Estimation of Network Games Management Science, 2021, vol. 67, no. 12, pp. 7575–7598, © 2021 The Author(s)

**Proposition 1** (Local Stability Condition). Under the condition that  $|\gamma_i| \cdot deg_i \cdot f^{\max}(\cdot) < 1$ , where  $deg_i$  is the weighted degree of consumer *i* (*i.e.*,  $deg_i = \sum_{j \in N_i} a_{ij}$ ) and  $f^{\max}(\cdot)$ , the maximal value of  $f(\cdot)$ , we have (1) the equilibrium is locally stable; and (2) the locally stale equilibrium is also unique.

*Proof:* See Online Appendix A for the proof.

Proposition 1 implies that if the susceptibility to social-influence parameters are in parameter space for all consumers, with  $|\gamma_i| < 1/(deg_i \cdot f^{\max}(\cdot))$ , we obtain a locally stable equilibrium that is also unique. Note that this sufficient condition must hold for all individuals and extends the unique equilibrium restriction with homogenous parameters (Lee et al. 2014) to the case with individual-specific parameters.<sup>4</sup> Allowing for individual-specific parameters results in individual-specific restrictions that are less stringent compared with the homogeneous situation. The reason is as follows: In the case of homogeneous parameters, the restriction for all consumers is determined by the inverse of the maximum weighted degree in the network; in contrast, in the case of individualspecific parameters, the restriction is individual-specific and only depends on a consumer's own weighted degree. Moreover, we also show that the unique equilibrium is locally stable, which implies that the equilibrium is a likely outcome of the game. Finally, the uniqueness of the locally stable equilibrium enables us to apply the stochastic Bayesian Markov  $chain\,Monte\,Carlo\,(MCMC)\,algorithm.\,The\,algorithm$ significantly reduces computational complexity and enables us to estimate the proposed social-influence model on large social networks.

To further understand the local stability condition, consider a two-node social network with two connected consumers, 1 and 2. An intuitive process to understand the choice equilibrium of the proposed model is to assume an initial probability of consumer 1 choosing  $d_1 = 1$ . This consumer influences the other consumer, which changes the choice probability of consumer 2. Subsequently, because consumer 2 changes its choice probability and influences consumer 1 as well, consumer 1 changes its own choice probability. This procedure repeats until both consumers reach a steady state, which equals the equilibrium solution of the model. The susceptibility to social-influence parameter serves as a factor that discounts the order of influence (e.g., first order: consumer 1 influences consumer 2; second order: consumer 2 influences consumer 1 after being influenced by consumer 1; etc.). The restriction on the susceptibility to social-influence parameters ensures that the system converges to a unique steady state given the (best-response) dynamic process described above. Conceptually, the restriction implies that social influence decreases as it propagates over the social network. This property is especially appealing if one adds a third consumer to the example that is only connected to consumer 2. The change of behavior of consumer 2 due to the change of behavior of consumer 1 will affect consumer 3. However, this indirect effect of consumer 1 on consumer 3 is weaker than the direct effect of consumer 1 on 2, because of the restriction on social influence. Similar diminishing effects will occur if the network is extended with more consumers, wherein the social effects become weaker as the distances between two consumers increase.

Finally, although previous research derived the explicit form of the unique equilibrium when the error distribution F is a type-I extreme value distribution (Brock and Durlauf 2001, Lee et al. 2014), we derive it for a general form of F in the following proposition.

**Proposition 2** (Explicit Form of the Unique Locally Stable Equilibrium). *Given the unique locally stable condition (Proposition 1), the equilibrium of this game defined by Equation (5) can be expressed as an infinite composition function evaluated at the intrinsic choice probability*  $F(\Delta)$ :

$$P^* = G_{\infty}(F(\mathbf{\Delta})|\mathbf{\Omega})$$

 $G_{\infty}(z|\Omega) = \lim_{m \to \infty} G_m(z|\Omega)$  is the infinite composition function of  $G(z|\Theta)$ , with  $G_m(z|\Omega) = G \circ G \circ ... \circ G(z|\Omega)$ 

and  $G(z|\Omega) = F(\Delta + \Gamma Az)$ , where  $\Omega = \{\Delta, \Gamma, A\}$ .

*Proof:* See Online Appendix A for the proof.

Proposition 2 provides an exact formula for the equilibrium choice probability, which equals  $G_{\infty}(F(\Delta)|\Omega)$ . This formula is particularly useful to express the intuition behind our counterfactual analysis, as we show in Sections 4.4.1 and 4.5.

## **3. Identification and Estimation Method 3.1. Model Identification**

In line with previous research of static games (Bajari et al. 2010, Ellickson and Misra 2011, Vitorino 2012) and dynamic discrete choice models (Rust 1987, Hotz and Miller 1993, Magnac and Thesmar 2002), we assume i.i.d. error terms  $\varepsilon_{it}$  across consumers and time (Section 2.1). In practice, the i.i.d. assumption may be challenging, as friendships are not formed randomly (i.e., homophily), and consumers may be exposed to the same external shocks (i.e., exogenous effects), resulting in a possible omitted variable bias (Manski 1993). This challenges the identification of the susceptibility to social-influence parameter  $\gamma_i$ . To address the challenge, we assume that the researcher observes T independent realizations of the game. This allows us to estimate individual fixed effects and consider individual-specific parameters for other variables in Equation (1) with a random-coefficient specification. The inclusion of individual fixed effects controls for homophily by capturing unobserved characteristics that may be correlated across consumers (Narayanan and Nair 2013, Shriver et al. 2013).<sup>5</sup> In addition, to control for external shocks that are common to all consumers, we include time fixed effects.

However, the inclusion of common exogenous shocks may not fully address the concerns over identification, because there may be unobserved variables that temporarily affect a subset of connected consumers. If omitted, these variables may be absorbed by the socialinfluence parameters, resulting in upward-biased estimates of social-influence parameters. In such a scenario, the identification of social-influence parameters hinges on the inclusion of excluded variables, which directly shift the utilities of some consumers and indirectly affect the utilities of others through social influence (see Bajari et al. 2010 and Vitorino 2012 for more details). Previous research often used group-specific shocks as excluded variables, such as those of school grades, families, or neighborhoods (De Giorgi et al. 2010, Lin 2010, Nicoletti et al. 2018). This is because consumers in the same group are more likely to be exposed to the same external shocks. For example, Shriver et al. (2013) grouped surfers into geo-locations and used varying wind speeds across days at these locations as instruments for surfers' blogging activities. However, in many applications, such time-varying excluded variables are difficult or even impossible to obtain (Imbens 2014), which limits the application of the proposed structural model. To address this challenge, we propose to exploit the network structure to identify subsets of consumers who are likely to be members of the same group and, hence, affected by the same groupspecific common temporal shocks. To this end, we utilize community-detection algorithms (Newman 2006, Blondel et al. 2008, Rosvall and Bergstrom 2008, Ronhovde and Nussinov 2009).

Previous research shows that social networks consist of different communities, which tend to have a denser network structure (Fortunato 2010). Consumers from the same community often have similar interests, speak the same languages, live in the same geographic regions, and bear other similarities. An important advantage of community detection is that latent communities are often more predictive of actual social group memberships than individual characteristics (Yang et al. 2013). For example, Fortunato (2010) found that results from communitydetection algorithms accurately predict how members of a karate club were divided into two social groups after a conflict happened between members. Such prediction was difficult to produce with individual characteristics alone. Similarly, Yang et al. (2013) found that network connections are more predictive of actual social group membership across multiple social media sites (e.g., Flickr, Twitter,

Google+, and Facebook) than clusters formed based on thousands of individual characteristics. Hence, similar to those in observed groups (e.g., school grades and neighborhoods), people in different communities are likely to be exposed to different "localized" external temporal shocks. To capture the effects of these external shocks, we propose to use communityspecific time fixed effects as excluded variables that only affect people in the same community at a certain time, but are excluded from the decisions of others in different communities.

Our community-detection approach to identify social influence is particularly attractive when observed excluded variables are unavailable to researchers. However, even if such information is available, the excluded variables constructed from communities can complement the observed ones (Newman and Clauset 2016). Moreover, based on the findings of Yang et al. (2013), it is plausible that excluded variables based on community-specific time fixed effects are superior to time fixed effects of observed groups, as communities are often more related to behaviors of consumers.<sup>6</sup>

#### 3.2. Challenges in Estimation

Previous literature suggested two approaches to estimate discrete choice models with interactions. The first approach consists of an iterative process with two steps (Zhu and Singh 2009, Bajari et al. 2010). The first step solves the system of nonlinear equations (e.g., Equation (5)), and the second step maximizes the likelihood function given the equilibrium derived in the first step. The second approach, named Mathematical Programming with Equilibrium Constraints, involves a direct and exhaustive search over the parameters and equilibria space (Su and Judd 2012). Such a constrained optimization approach has been applied in marketing to a static entry game with incomplete information (Vitorino 2012). A common limitation of the two approaches is the high computational complexity. Previous research only applied these methods to empirical settings with groups that are substantially smaller than online social networks. For example, Zhu and Singh (2009) considered the entry decisions of three retailers and Vitorino (2012) up to nine stores in each market. Lee et al. (2014) adopted the first estimation approach and applied their model to a collection of networks, each with at most a few hundred people. However, their estimation method is infeasible for large social networks with heterogeneous consumers. In sum, it is infeasible to apply previous estimation procedures to

empirical settings of large social networks in the current research.

An important reason for the complexity of previous estimation procedures is the consideration of multiple equilibria.<sup>7</sup> Our modeling context admits complex interactions (both strategic complementarities  $\gamma_i > 0$ and substitutions  $\gamma_i < 0$ ) between many consumers with individual-specific parameters. To reduce the complexity of the estimation, we ex ante select a unique equilibrium based on the local stability condition (see Proposition 1), following the theory of learning in games (Fudenberg and Levine 1998). Such a treatment is important for our research, as it enables us to apply a stochastic Bayesian MCMC estimation procedure with the insights from Imai et al. (2009), which significantly reduces computational complexity. The procedure was originally developed to estimate dynamic discrete choice models and has been applied in structural demand estimation by Sun and Ishihara (2018). Similar to our model, dynamic discrete choice models contain a value function, which is the unique solution to a nonlinear system of equations. Subsequently, the value function is introduced into the likelihood function for parameter estimation. The stochastic Bayesian MCMC procedure eases the need to solve the nonlinear system of equations. The procedure derives a "pseudo" solution, instead of an exact solution that is computationally demanding. By doing so, the procedure replaces a complex loop with a direct evaluation step, which significantly reduces computational demands. Another advantage of the stochastic estimation procedure is that it allows for the natural incorporation of individual-specific parameters through a hierarchical Bayesian structure.

#### 3.3. The Stochastic MCMC Algorithm

Next, we present a detailed description of the stochastic Bayesian MCMC algorithm.

The Stochastic Bayesian MCMC Algorithm. Start with initializing a guess of the solution to the system of equations and a parameter vector—that is,  $\langle \hat{p}_i^{*(0)}, \Theta_i^{(0)} \rangle$ . At iteration (*r*), we execute the following steps:

Step 1: Use pseudo-solution  $\hat{p}_i^{*(r-1)}$  in the likelihood function to draw parameter vector  $\Theta_i^{(r)}$  using a standard MCMC procedure (e.g., Gibbs-sampling).

Step 2(a): Derive a candidate pseudo-solution  $\tilde{p}_i^{*(r)}$ based on the history of iterations stored in memory<sup>8</sup>  $H_i^r = \{\hat{p}_i^{*(l)}, \Theta_i^{(l)}\}_{l=r-R'}^{l-1}$ , where the term  $\tilde{p}_i^{*(l)}$  is the pseudosolution at iteration (*l*),  $\Theta_i^{(l)}$  is the parameter vector drawn at iteration (*l*) and *R* is set by the researcher and represents the length of the history (i.e., how many previous iterations). The candidate pseudo-solution  $\tilde{p}_i^{*(r)}$  is a weighted average of previous pseudo-solutions, which is computed as follows:

$$\tilde{p}_{i}^{*(r)}(H_{i}^{r}) = \sum_{l=r-R}^{r-1} \omega \left(\Theta_{i}^{(l)}, \Theta_{i}^{(r)}\right) \hat{p}_{i}^{*(l)}, where$$

$$\omega \left(\Theta_{i}^{(l)}, \Theta_{i}^{(r)}\right) = \frac{K_{h}\left(\Theta_{i}^{(l)}, \Theta_{i}^{(r)}\right)}{\sum_{l=r-R}^{r-1} K_{h}\left(\Theta_{i}^{(l)}, \Theta_{i}^{(r)}\right)}.$$
(7)

In the above equation,  $K_h(\cdot)$  is a multivariate kernel density with bandwidth h. In theory, any form of density function can be applied here (e.g., a Gaussian kernel).

Step 2(b): Using the candidate pseudo-solution  $\tilde{p}_i^{*(r)}$ , we apply the operation defined in Equation (5) only once to obtain the final pseudo-solution, with  $\hat{p}_i^{*(r)} = F(V_i + \gamma_i \sum_{j=1}^{N_i} a_{ij} \tilde{p}_j^{*(r)})$ . Store this value  $\hat{p}_i^{*(r)}$  along with parameter vector  $\Theta_i^{(r)}$  in memory and go to iteration (r + 1).

Our algorithm differs from Imai et al. (2009) in two aspects. First, all the posterior distributions in Step 1 are standard, and, thus, we can use a Gibbs sampler, which generally converges faster than Metropolis– Hastings steps. Second, Steps 2(a) and 2(b) use policyfunction iterations to approximate the equilibrium solution of the game (Equation (5)), instead of valuefunction iterations to approximate the fixed point of the Bellman equation.

To prove that the parameter estimates converge to the true posterior distribution, we show that Steps 2(a) and 2(b) converge to the true equilibrium of the game defined in Equation (5).<sup>9</sup> Suppose the parameter sets  $\Theta_i^{(r)}(\forall i = 1, \dots, N)$  stay fixed at a value  $\Theta_i^*$  for each iteration r. With fixed parameters, weights  $\omega$  in Equation (7) for all historical draws equal 1/R. Assume that we start with a vector of choice probabilities  $p^0$ . In the first iteration, we derive  $p^1 = F(\Delta + \Gamma A p^0)$ , and, subsequently, in the second iteration,  $p^2 = F\Delta + \Gamma A \times$  $(\frac{1}{2}p^0 + \frac{1}{2}p^1)$ . Given the contraction-mapping property (see proof for Proposition 1), we derive the following,

$$\left\| \boldsymbol{p}^{2} - \boldsymbol{p}^{1} \right\|_{\infty} = \left\| F \left( \boldsymbol{\Delta} + \Gamma A \left( \frac{1}{2} \boldsymbol{p}^{0} + \frac{1}{2} \boldsymbol{p}^{1} \right) \right) - F \left( \boldsymbol{\Delta} + \Gamma A \boldsymbol{p}^{2} \right) \right\|_{\infty}$$
$$\leq \frac{1}{2} \lambda \left\| \boldsymbol{p}^{1} - \boldsymbol{p}^{0} \right\|_{\infty}, \tag{8}$$

where  $||u - v||_{\infty} = \max(|u - v|)$  is the maximum distance between two elements of vectors u and v and  $\lambda \in (0,1)$ . The same relationship holds for additional iterations. The contraction-mapping property ensures that the discrepancy between derived choice probabilities decreases across iterations, and the stochastic algorithm eventually converges to the equilibrium of the game. As noted by Imai et al. (2009), the stochastic algorithm gives more weights to equilibrium-choice

probabilities that are closer to the current parameter draw and ensures the convergence when parameters vary across iterations.

In sum, the proposed stochastic algorithm guarantees that the parameter simulations converge to the true posterior distributions after a sufficiently large number of iterations. Online Appendix B.2 further provides a detailed description of the full stochastic Bayesian MCMC procedure that we applied in the empirical application. Moreover, simulation studies reveal that the proposed method recovers the parameters accurately, while significantly reducing computational complexity, as shown next.

#### 3.4. Scalability Analysis of the Algorithm

To examine the performance of the stochastic Bayesian MCMC procedure, we demonstrate parameter recovery, as well as CPU time for model estimation. All models were estimated on a standard desktop with an Intel<sup>®</sup> Xeon<sup>™</sup> E-2176M Processor (six core, 2.70 GHz, 4.40-GHz Turbo, 12-MB cache) and 32-GB DDR4 internal memory. We compared the proposed estimation procedure (i.e. "Stochastic") with the traditional Bayesian MCMC method (i.e. "Deterministic"). In contrast to the stochastic Bayesian MCMC procedure, the deterministic method follows a full iterative approach to solve the equilibrium-choice probabilities, which is similar to maximum-likelihood approaches used in previous literature.

**3.4.1. Simulating Data.** To generate the data, we used the following specification for the utility function (see Equation (1)):

$$U(d_{it} = 1 | x_{it}, d_{-it}, \Theta_i)$$
  
=  $\beta_{0i} + \lambda_t + \beta_{1i} x_{it} + \lambda_t^{\text{com}} x_i^{\text{com}} + \gamma_i \sum_{j \in N_i} \tilde{a}_{ij} p_{jt}^* + \varepsilon_{it}.$ 
(9)

In Equation (9), we assume a decision-making process of binary choices  $d_{it}$  of consumer *i* at period *t* and random shocks following an i.i.d. standard normal distribution with  $\varepsilon_{it} \sim N(0,1)$ . We row-normalize the adjacency matrix as in the empirical application, with  $\tilde{a}_{ij}$  as the *ij*th entry of the row-normalized adjacency matrix *A* and  $\sum_{j \in N_i} \tilde{a}_{ij} = 1$ . With the row-normalized adjacency matrix  $\tilde{A}$  and standard normal errors  $\varepsilon_{it}$ , the susceptibility to social-influence parameters  $\gamma_i \in$  $(-\sqrt{2\pi},\sqrt{2\pi})$ , a direct result from applying the local stability condition in Proposition 1. Individual and time fixed effects are captured by  $\beta_{0i}$  and  $\lambda_t$ , respectively. We allow choices to be affected by an individual- and time-varying explanatory variable  $x_{it}$  that is simulated from a standard normal distribution. Finally, we divide consumers into two communities of equal sizes and then create a community dummy  $x_i^{\text{com}}$  that equals one if individual *i* is from the first community and zero otherwise, as the second community serves as baseline. Finally, community-specific time fixed effects are captured by  $\lambda_t^{\text{com}}$ .

In our simulation analyses, we considered both situations with and without individual-specific parameters. In the homogeneous case, we set  $\beta_{1i}$  and  $\gamma_i$  to the same values for all consumers. In the heterogeneous case, we generate individual-specific parameters from normal distributions common to all consumers, with  $\beta_{1i} \sim N(\mu_{\beta_i}, \sigma_{\beta_i}^2)$  and  $\gamma_i \sim N(\mu_{\gamma}, \sigma_{\gamma}^2)I(-\sqrt{2\pi} < \gamma_i < \sqrt{2\pi})$ . In both cases, we simulated individual-fixed effects, time-fixed effects, and community-specific effects from uniform distributions as specified below. Using these variables and parameters, we computed equilibrium choice probabilities by solving the system of Equations (5) with the convergence tolerance set to  $10^{-15}$ . After computing the equilibrium choice probabilities, we simulated latent utilities and derived users' binary choices accordingly.

Given the setup, we investigate the computational demands as a function of sample size and the magnitude of the susceptibility to social influence. Sample size affects the computational demand of each iteration, and the magnitude of social influence affects the rate of convergence (i.e., how many iterations are needed) of the equilibrium choice probabilities. In our simulations, we generated 10 different sample sizes, with *N* set to {200, 400,...,2,000}. The number of time periods T was fixed to 100, similar to the empirical application. As we included community-specific time fixed effects, we followed Girvan and Newman (2002) to simulate networks with known communities.<sup>10</sup> The density of the simulated networks was kept at 0.01, which results in, on average, 2,4,...,20 friends for different sample sizes. We simulated four scenarios for both the homogeneous and heterogeneous parameter cases, with the true values of susceptibility to social-influence parameters set to, respectively,  $\{-1.00, -0.01, 0.01, 1.00\}$ . In the heterogeneous cases, we set the corresponding variance  $\sigma_{\nu}^2$  to {0.25, 0.01, 0.01, 0.25}. In all scenarios, we simulated individual fixed effects  $\beta_{0i}$  from a uniform distribution between -1.00and 1.50. The mean of parameter  $\beta_{1i}$  was set to one, with a variance  $\sigma_{\beta_1}^2 = .25$  in the heterogeneous case. The time fixed effects  $\lambda_t$  and community-specific time fixed effects  $\lambda_t^{\text{com}}$  were simulated from a uniform distribution between -0.50 and 0.50. The combination of all factors leads to 40 scenarios for either the homogeneous or the heterogeneous case (10 sample sizes  $\times$  4 social-influence conditions).

**3.4.2. Simulation Results: Comparison Between the Deterministic and Stochastic Algorithm.** In the estimation, we ran a total of 10,000 iterations and used the

first 5,000 draws as the burn-in period, long after the convergence of the runs. For the stochastic method, we set the length of history R to 20 and chose a Gaussian kernel with bandwidth equal to the rule-of-thumb procedure proposed by Scott and Sain (2005). For all estimation runs, we set the initial values of all parameters to half of their true values. The first two column groups of Table 1 present the true and estimated values of social-influence parameters with the deterministic and stochastic estimation procedure.<sup>11</sup> Estimation results of both algorithms are virtually identical, and all 95% posterior intervals contain the true parameter values.

As discussed above, the main advantage of the stochastic method is the improvement of the efficiency of estimation. To examine the efficiency gain, we recorded the CPU time of each iteration for all simulation conditions. As the computational demand depends on the magnitude of social influence, but not the sign, we combined the simulation results with the same absolute values of (mean) social influence. Figure 1 summarizes the computational demands (seconds per 10 iterations) of both methods across different scenarios. First, in all scenarios, the stochastic method shows substantial efficiency gains over the deterministic method, especially for larger sample sizes and social-influence parameters of larger magnitude. Second, although the computational complexity of the deterministic method is polynomial in the sample sizes, the complexity of the stochastic method is linear, which makes it especially valuable for large sample sizes. Third, the computational complexity of the stochastic method does not depend on the magnitude of social influence, whereas the deterministic method is substantially more demanding for social influence of larger magnitude. In conclusion, the proposed stochastic estimation procedure accurately recovers true parameters, but significantly reduces the computational burden.

**3.4.3. Simulation Results: The Importance of Communities.** To investigate the importance of community-specific time fixed effects for the identification of social-influence parameters, we further estimated models that ignore these effects (i.e.,  $\lambda_t^{\text{com}} = 0$  in Equation (9)) on the same data simulated in Section 3.4.1. The last column group of Table 1 presents the estimation results of social-influence parameters and Online Appendix C.2 reports the estimates of the fixed effects. In line with the expectation, the mean estimates of social influence are biased upward, especially in situations with stronger social influence. For both the homogenous and heterogenous cases, the 95% posterior intervals do not contain the true values of social influence. These results show that ignoring

	Deterministic approach		Stochastic	Stochastic approach		Without communities	
True parameter	<i>N</i> = 200	<i>N</i> = 2,000	<i>N</i> = 200	<i>N</i> = 2,000	<i>N</i> = 200	N = 2,000	
Homogeneous: $\gamma$							
-1.00	-0.98 (-1.10, -0.86)	-0.99 (-1.11, -0.87)	-0.98 (-1.09, -0.85)	-1.00 (-1.09, -0.89)	-0.87 (-0.99, -0.74)	-0.88 (-1.00, -0.76)	
-0.10	-0.11 (-0.23, -0.01)	-0.10 (-0.18, -0.02)	-0.11 (-0.20, -0.01)	-0.09 (-0.18, -0.01)	-0.08 (-0.20, 0.03)	-0.09 (-0.23, 0.08)	
0.10	0.10 (0.00, 0.22)	0.10 (0.05, 0.20)	0.10 (0.01, 0.19)	0.10 (0.02, 0.19)	0.13 (0.01, 0.25)	0.12 (-0.02, 0.25)	
1.00	1.04 (0.89, 1.18)	1.01 (0.80, 1.17)	1.04 (0.91, 1.18)	1.02 (0.88, 1.14)	1.19 (1.04, 1.32)	1.19 (1.05, 1.36)	
Heterogeneous: $\mu_{\gamma}$							
-1.00	-1.03 (-1.18, -0.88)	-1.00 (-1.13, -0.87)	-1.02 (-1.18, -0.87)	-1.01 (-1.13, -0.90)	-0.82 (-0.97, -0.66)	-0.83 (-0.97, -0.69)	
-0.10	-0.09 (-0.22, 0.04)	-0.10 (-0.15, -0.08)	-0.10 (-0.26, 0.04)	-0.09 (-0.17, -0.04)	-0.08 (-0.20, 0.07)	-0.08 (-0.17, 0.01)	
0.10	0.09 (-0.04, 0.21)	0.09 (0.03, 0.15)	0.09 (-0.01, 0.19)	0.09 (0.01, 0.14)	0.12 (-0.01, 0.25)	0.12 (0.06, 0.19)	
1.00	0.98 (0.84, 1.15)	0.99 (0.85, 1.13)	0.97 (0.83, 1.12)	0.98 (0.84, 1.13)	1.19 (1.05, 1.34)	1.20 (1.03, 1.37)	
Heterogeneous: $\sigma_{y}^{2}$							
0.25	0.25 (0.07, 0.48)	0.25 (0.09, 0.41)	0.26 (0.07, 0.50)	0.26 (0.10, 0.42)	0.29 (0.11, 0.53)	0.30 (0.12, 0.49)	
0.01	0.01 (0.00, 0.06)	0.01 (0.00, 0.04)	0.01 (0.00, 0.05)	0.01 (0.00, 0.04)	0.01 (0.01, 0.06)	0.01 (0.00, 0.07)	
0.01	0.01 (0.00, 0.06)	0.01 (0.00, 0.04)	0.01 (0.00, 0.05)	0.01 (0.00, 0.04)	0.01 (0.00, 0.10)	0.01 (0.00, 0.04)	
0.25	0.26 (0.12, 0.49)	0.26 (0.10, 0.41)	0.26 (0.10, 0.43)	0.25 (0.12, 0.41)	0.20 (0.16, 0.36)	0.21 (0.07, 0.37)	

Table 1. Estimates of Social-Influence Parameters in Simulation Studies

Note. The 95% posterior confidence intervals are in the parentheses.

community-specific time fixed effects results in substantial biases, even if individual fixed effects are considered. This is because individual fixed effects control for time-invariant unobservables, but not time-varying unobservables.<sup>12</sup> Overall, the simulation results highlight the importance of considering communities for the identification of social-influence parameters.



← Deterministic.(.10) -▲- Stochastic.(.10) -■ Deterministic.(1.00) +- Stochastic.(1.00)



# 4. Empirical Application: Login Decisions in an Online Game

To illustrate the applicability of our model, we obtained a data set that contains information about Asian users from a massive multiplayer online roleplaying game. Online gaming has become a multibilliondollar industry with the global revenue of \$83.10 billion in 2019.<sup>13</sup> Importantly, online games provide useful data about social and economic interactions (Bainbridge 2007) and are thereby ideal empirical contexts to study social influence. The online game we study is based on a Western fairytale storyline and was globally one of the largest games at the time of the data collection. Like many other online games, the game is free to play and relies on in-game purchases of users as the revenue source (i.e., a freemium model). Therefore, it is important for the online game to keep users active by stimulating users' logins to the game. In the game, users form friendships, and social influence is an important driving factor of user behaviors. Therefore, it is important to keep track of those influential users, whose login activities stimulate more frequent logins of other users. We observe login decisions to the game (as 0–1 decisions) and users' profiles, such as gender and geo-locations (i.e., cities), as well as network connections in the game.

Login decisions are likely to follow our modeling assumptions, as the utility of a login for a user depends on the online status of other users. For example, users may share information and complete tasks together (e.g., slaving monsters and finding treasures). However, it is difficult for users to coordinate or observe the login decisions of their friends before logging in online. Moreover, even if users can communicate outside the game and inform their friends that they are likely to play during a certain time period, friends may not know exactly when they would be online. Therefore, users may respond to the expectations of the logins of their friends, which depend on the time of login decisions and characteristics of their friends. Following the assumptions in Section 2, users are assumed to be myopic, and their login decisions in different time periods are, therefore, independent realizations of the same network game. To provide some support for this assumption, we followed a procedure similar to Gruber and Köszegi (2001), which used a preannouncement of an increase in the cigarette tax to test the forwardlooking tendencies of smokers. In our data, we observe an in-game preannouncement of a major game update. The announcement happened on Day 26 and lasted until the end of the observation period, which was just before the actual update. The update made available new content (i.e., new tasks and territories), and some of them could only be accessed if users

reached a minimum level. Hence, if users were forwardlooking, we would expect an increase in their logins during the announcement period, as increased logins would enable them to reach the minimum levels required by the new contents. We ran a probit regression of login decisions with all exogenous variables that we used in the full model, such as individual and time fixed effects, and the announcement dummy. The results revealed no empirical evidence of forwardlooking tendencies (the coefficient of the announcement dummy is -0.014, with standard error 0.021).<sup>14</sup>

#### 4.1. Data Description

The data contain 30 days of login records of 25,418 users in a complete social network of the online game. To start gaming, users must select one game server, choose a character from a range of hero classes (e.g., warriors, archers), and use their selected avatars to enjoy the game contents. Given that many in-game tasks require team efforts, users have incentives to form in-game friendships. Notably, users can only befriend others on the same game server. We thus observed a full network of all users (around 110,000 users) from one server. Given the full network, we obtained a giant component of around 30,000 users. In this giant component, any two users can reach each other, and none of the users has friends outside of the giant component. We then removed inactive users, who did not log in during the observation period of 30 days, to obtain the final data set of 25,418 active users.

Although users can always change their friendships in the game, we observed a mature and stable social network with only few changes. During the observation period, only 130 users (0.051%) formed new friendships with one another (i.e., 65 new pairs of friends), and none of them defriended any other users. New friendships were formed on 24 out of 30 days (i.e., days without any new friendships are {6, 13, 14, 18, 21, 26}), and we did not find any upward or downward trends in friendship formation (Man-Kendall test: Z = 0.66, p = 0.51).<sup>15</sup> At the end of the observation period, the average degree of the network is 5.87 (with standard deviation: 6.98). The degree distribution (d) is power-law like with  $p(d) = (\hat{\alpha} - 1)d^{-\hat{\alpha}}$ , where parameter  $\hat{\alpha}$  is estimated at 1.17 (standard error: 0.001). The mean percentage of common connections between two friends is 2.57%, with a minimum overlap of 0.00%, median 0.00%, and maximum 50.27% (see Structural Equivalence in Table 2).

We observed login activities of users over a period of 30 days. To capture login decisions, we sliced days in quarters (i.e., 12 a.m.–6 a.m., 6 a.m.–12 p.m., 12 p.m.–6 p.m., and 6 p.m.–12 a.m.) and constructed a dummy variable that equals one if a user logs into the game during a quarter of the day, and zero otherwise.

Variables		Mean	Standard deviation	Min	Max	Correlation (with $Y$ )
Login	Y	0.16	0.37	0.00	1.00	1.000
Social Influence	Average Friends' Logins	0.21	0.29	0.00	1.00	0.212
Exogenous Effects	Regular Promotion	0.04	0.20	0.00	1.00	0.002
Targeted Promotions	Promotion (Women's Day)	0.02	0.13	0.00	1.00	0.004
	Promotion (White Day)	0.02	0.13	0.00	1.00	0.006
Network Measures	Degree Centrality	5.87	6.98	1.00	76.00	0.13
	Structural Equivalence	0.03	0.05	0.00	0.50	0.02

**Table 2.** Descriptive Statistics of Variables

Because of the inclusion of individual fixed effects, we did not incorporate time-invariant individual characteristics, such as gender (48.76% are female) and geographical location (i.e., users were in 211 cities). Importantly, we observed time-varying in-game promotions, of which one was a regular promotion for all users and two were targeted to users of either gender. The regular promotion rewarded users with in-game items (e.g., magic potions and costumes) if they logged in during the promotion. The two targeted promotions corresponded to popular public events, International Women's Day and White Day,<sup>16</sup> which were targeted at females and males, respectively. Similar to the regular promotion, during the targeted promotions, targeted users were rewarded with valuable in-game items if they logged in. Hence, the targeted promotions were exogenously determined and shifted logins of only a selected group of users. Moreover, most users befriend both females and males and have on average 48.98% (standard deviation (std.): 33.69%) of friends from the opposite sex (females: 45.03% with std. 31.63% and males: 52.73% with std. 35.13%). The targeted promotions are, thus, valuable sources of variation that we can exploit for the identification of social influence. For the regular promotion, we used a dummy with one denoting the quarter-days with the promotion, and zero otherwise. For the targeted promotions, we used two dummy variables with one denoting a user targeted during the quarter-days of the promotions, and zero otherwise.

Summary statistics of all variables are presented in Table 2. As the first evidence of social influence, we find a positive correlation between users' own logins and the average logins of direct friends (r = 0.21; on average, 35.31% of friends are online when a user logs in and 18.57% when a user does not log in). Moreover, we also find significant, but weaker, correlations between users' own logins and the average login decisions of their second-order (r = 0.19) and third-order (r = 0.13) friends. The differences between the correlations are significant (first versus second order: Z score = 27.59, p < 0.00; second versus third order: Z score = 108.81, p < 0.00), which indicate that social

influence may weaken over the network. Such a pattern is in line with the implications of the unique equilibrium condition. As expected, all promotions show positive and significant (p < 0.00) correlations with logins. Finally, we also explored whether login decisions varied across time (see Figure 2) and found that users were less likely to login during the first quarter (12 a.m.–6 a.m., average login probability 6.67%) and second quarter (6 a.m.–12 p.m., average login probability 10.92%), compared with the third quarter (12 p.m.–6 p.m., average login probability 22.97%) and fourth quarter (6 p.m.–12 a.m., average login probability 24.08%).

#### 4.2. Community Detection

As discussed in Section 3.1, we propose to use community-specific time fixed effects to facilitate the identification of social influence. To do so, we must divide the social network into different communities, with many connections within communities and few connections between communities. We followed Lancichinetti and Fortunato (2009), who recommended four types of algorithms for undirected and unweighted networks. These algorithms are as follows: (1) Spectral Bisection (Newman 2006), (2) Infomap (Rosvall and Bergstrom 2008), (3) Blondel (Blondel et al. 2008), and (4) RN (Ronhovde and Nussinov 2009). For the first three methods, we used the C++ igraph library (Csardi and Nepusz 2006). We used the RN algorithm implemented in C++ (Ronhovde and Nussinov 2009) and specified 11 different parameter values to weight connections within and between communities (i.e., the  $\gamma$  parameter in Ronhovde and Nussinov 2009, which we set to 0.01, 0.1, 0.2, 0.3... and 1). In total, we compared 14 community-detection algorithms. To select the optimal solution, we adapted four quality measures based on Leskovec et al. (2010). Intuitively, a good community-detection solution should result in communities where people are densely linked within communities and sparsely between communities. The adapted measures reflect this fundamental criterion. Specifically, the four quality measures are as follows:



#### Figure 2. Aggregate Logins over Time

(1) Conductance (i.e., the percentage of connections between communities over all the connections in the network), (2) Expansion (i.e., the average number of friends an individual has outside of his community), (3) Average Out-Degree Friends (Average-ODF; i.e., the average percentage of friends that an individual has outside of his/her community), and (4) Flake-ODF (i.e., the percentage of people that have fewer friends within than outside of their own communities). For each criterion, lower values indicate better performances. Table 3 reports the performance measures for each of the 14 community-detection algorithms. Overall, the Blondel algorithm provides the best performance on all measures, except for the Flake-ODF measure. It has a similar, but slightly lower, performance compared with RN algorithm with parameter 0.01 (0.0044 versus 0.0001, respectively, for Blondel and RN algorithm). Therefore, we decided to use the communities detected by the Blondel algorithm in our follow-up analysis.<sup>17</sup>

Using the algorithm, we detected 125 communities, with sizes ranging from 26 to 817 users.<sup>18</sup> For an average user, 80.7% of his friends are in the same community (min: 7.7%, median 66.7%, max: 100%) and, on average, belong to 3.12 different communities (min: 1.00, median: 2.00, max: 24.00). This variety in community-membership of friends facilitates the identification of social influence, on the condition that different communities experience different external shocks. To test this, we ran an analysis of variance (ANOVA) on the average logins of communities across time, with communities and various time fixed effects as factors. In particular, we included quarter-ofthe-day effects, day-of-the-week effects, and linear and quadratic trends. The ANOVA results reveal the significant main effect of communities (F = 20.75, p < 0.00), as well as significant interactions between communities and quarter-of-the-day effects (F = 20.65, p < 0.00), linear trend (F = 40.76, p < 0.00), and quadratic trend

Table 3. Performance of Different Community-Detection Algorithms

Algorithms	Conductance	Expansion	Average-ODF	Flake-ODF
Bisection	0.3700	2.1735	0.2547	0.1371
Infomap	0.4608	2.7064	0.2963	0.2220
Blondel	0.0032*	0.0190*	0.0062*	0.0044
RN $(r = 1.0)$	0.3955	2.3232	0.2705	0.1551
RN $(r = 0.9)$	0.3722	2.1864	0.2568	0.1332
RN $(r = 0.8)$	0.3412	2.0043	0.2393	0.1083
RN $(r = 0.7)$	0.2998	1.7608	0.2128	0.0688
RN $(r = 0.6)$	0.2446	1.4369	0.1772	0.0426
RN $(r = 0.5)$	0.2012	1.1818	0.1459	0.0176
RN $(r = 0.4)$	0.0829	0.4868	0.0737	0.0168
RN $(r = 0.3)$	0.0292	0.1716	0.0314	0.0058
RN $(r = 0.2)$	0.0092	0.0539	0.0128	0.0042
RN $(r = 0.1)$	0.0041	0.0241	0.0066	0.0035
RN $(r = 0.01)$	0.0110	0.0647	0.0098	0.0001*

\*The best-performing algorithm on corresponding criterion.

(F = 5.82, p < 0.00). Figure 3 shows the average logins of different communities across different quarter days. As shown in the plot, the average logins vary substantially across quarter days, with some communities more active in the third quarter (12 p.m.–6 p.m.) and others in the fourth quarter (6 p.m.–12 a.m.). Overall, these variations make the inclusion of communityspecific time fixed effects a useful approach for the identification of social-influence parameters.

#### 4.3. Model Specification

For the empirical application, we specify the utility function (Equation (1)) as follows:

$$U(d_{it}=1|\mathbf{x}_{it}, d_{-it}, \Theta_i) = \beta_{0i} + \lambda_t + \beta_{1i} \mathbf{x}_t + \beta_2 \mathbf{x}_{it}^{\text{prom}} + \lambda_t^{\text{city}} \mathbf{x}_i^{\text{city}} + \lambda_t^{\text{com}} \mathbf{x}_i^{\text{com}} + \gamma_i \sum_{i \in N_i} \tilde{a}_{ij} p_{jt}^* + \varepsilon_{it}.$$

$$(10)$$

First, we assume a standard normal distribution for i.i.d. error term  $\varepsilon_{it}$ , which implies that login decisions follow a probit choice model.<sup>19</sup> As in the simulation studies, we row-normalized adjacency matrix A, so that all susceptibility to social-influence parameters are bounded with  $|\dot{\gamma}_i| < \sqrt{2\pi}$ .<sup>20</sup> We incorporated individual ( $\beta_{0i}$ ) and time fixed effects ( $\lambda_t$ ), and set the effects of two time periods (i.e., quarter day 1 and 85) to zero for identification.<sup>21</sup> The vector  $x_t$  contains the regular promotion. Finally, we included three sets of excluded variables based on the following variables: (1) two targeted promotions  $(x_{it}^{\text{prom}})$ , (2) a vector of 211 cities dummies ( $x_i^{\text{city}}$ ), and (3) a vector of 125 community dummies ( $x_i^{\text{com}}$ ). For the latter two, we set one city and one community as baselines. In the estimation, we allowed the effects of the regular promotion  $(\beta_{1i})$  and the susceptibility to social influence  $(\gamma_i)$  to be

individual-specific using random coefficient specifications with normal distributions. We did not allow the effects of targeted promotions  $\beta_2$  to be individual-specific, because only a subset of users were targeted. Finally, parameter vectors  $\lambda_t^{\text{city}}$  and  $\lambda_t^{\text{com}}$  capture, respectively, the city- and community-specific time fixed effects at time *t*.

#### 4.4. Estimation Results and Model Validation

To illustrate the importance of incorporating individualspecific parameters and examine the efficacy of different excluded variables, we compared our model to several nested alternatives. In addition, we also validated the identification strategy of incorporating community-specific time fixed effects by constructing exogeneity tests from a reduced-form analysis, as shown in Section 4.4.2.

4.4.1. Estimation Results and Model Comparisons. We estimated six different models, which all include individual fixed effects ( $\beta_{0i}$ ), exogenous effects ( $x_t$ ), and social influence ( $\gamma$ ). Model 6 is the full model as specified in Equation (10). Model 1 assumes that the effects of the regular promotion and social influence are homogeneous across users (i.e.,  $\beta_{1i}$  and  $\gamma_i$  are the same across all consumers). In contrast, all the remaining models allow for individual-specific parameters. Model 2 does not incorporate any excluded variables  $(\beta_2 = \lambda_t^{\text{city}} = \lambda_t^{\text{com}} = 0)$ , whereas model 3 only incorporates targeted promotions as excluded variables ( $\lambda_t^{\text{city}} = \lambda_t^{\text{com}} = 0$ ) and model 4 adds city-specific time fixed effects ( $\lambda_t^{\text{com}} = 0$ ). Finally, model 5 only includes community-specific time fixed effects ( $\beta_2$  =  $\lambda_t^{\text{city}} = 0$ ) and corresponds to situations where researchers do not directly observe any exogenous excluded variables. All models were estimated with

Figure 3. Average Logins of Different Quarters Across Different Communities Quarters of the Day • 00:00-06:00 • 06:00-12:00 • 12:00-18:00 + 18:00-24:00



the stochastic Bayesian MCMC method (see Online Appendix B.2 for the detailed MCMC procedure). We used 20,000 draws after a burn-in period of 10,000, long after the convergence (for the convergence diagnostics, see Online Appendix G).

Table 4 reports the estimation results of the models, as well as the fit statistics (Log Marginal Density (LMD)). The fit statistics indicate that the inclusion of individual-specific parameters strongly improves model performances (i.e., LMD of *Model*  $1 = -9.503 \times 10^{5}$ versus *Model*  $6 = -8.483 \times 10^{5}$ ). Moreover, compared with model 2, adding observed excluded variables (i.e., targeted promotions and city-specific time fixed effects) improves model fit (LMD of *Model*  $2 = -8.511 \times 10^5$  versus  $Model 3 = -8.504 \times 10^5$  versus  $Model 4 = -8.499 \times 10^5$ ). Interestingly, if we only add community-specific time fixed effects as in model 5, model fit improves much stronger, which indicates that communities may capture important differences between users (LMD of Model 5 =  $-8.495 \times 10^5$ ). Finally, the full model 6, which includes both observed excluded variables and community-specific time fixed effects, best describes the data (LMD *Model*  $6 = -8.483 \times 10^{5}$ ).

We next compare the estimates of social-influence parameters of different models. First, ignoring individualspecific parameters (model 1) significantly reduces the estimate of the susceptibility to social influence (mean estimate: 0.69 versus 0.91, respectively, for model 1 and model 6). Second, ignoring excluded variables significantly inflates the mean estimates of the susceptibility to social influence (1.11 versus 0.91, respectively, for model 2 and model 6, with all posterior draws of model 2 exceeding those of model 6). This highlights the importance of excluded variables for the identification of social influence. The parameter estimates of model 3 to model 5 further highlight the power of community-specific time fixed effects. Compared with model 2, the inclusion of targeted promotions and city-specific time fixed effects reduces the upward bias in the mean estimates of social influence (1.08 versus 1.03, respectively, for model 3 and model 4). However, these estimates are still significantly different from the full model (all posterior draws of models 3 and 4 are larger than those of model 6). In contrast, if only community-specific time fixed effects are included as excluded variables as in model 5, the mean estimate of social influence becomes close to that of the full model (mean estimates: 0.93 versus 0.91, respectively, for model 5 and model 6, with 4.63% of the posterior draws in model 6 larger than those of model 5). Figure 4 further illustrates the differences between the individual-specific estimates of the susceptibility to social influence for models 2-6. Model 5 produces a histogram of the

_	Model 1	Model 2	Model 3	Model 4	Model 5	Model 6
Variables	Homogeneous parameters	Without city and community and targeted promotions	Without city and community	Without community	Without city and targeted promotions	Full model
Social Influence						
Mean	<b>0.69</b> (0.65, 0.74)	<b>1.11</b> (1.08, 1.14)	<b>1.08</b> (1.05, 1.11)	<b>1.03</b> (1.00, 1.07)	<b>0.93</b> (0.89, 0.97)	<b>0.91</b> (0.87, 0.95)
Variance	—	<b>0.74</b> (0.70, 0.77)	<b>0.70</b> (0.68, 0.73)	<b>0.63</b> (0.61, 0.65)	<b>0.68</b> (0.66, 0.70)	<b>0.61</b> (0.59, 0.63)
Individual Fixed Effects	5					
Post hoc mean	-1.14	-1.71	-1.66	-1.56	-1.50	-1.47
Post hoc variance	0.48	0.42	0.40	0.34	0.33	0.38
Promotions						
Regular promotion	0.07	0.05	0.06	0.06	0.07	0.09
	(0.02, 0.12)	(0.04, 0.06)	(0.05, 0.08)	(0.00, 0.10)	(0.01, 0.13)	(0.02, 0.15)
Promotion (Women's Day)	<b>0.11</b> (0.09, 0.13)	_	<b>0.07</b> (0.05, 0.10)	<b>0.06</b> (0.04, 0.08)	_	<b>0.06</b> (0.04, 0.08)
Promotion (White Day)	<b>0.19</b> (0.17, 0.21)	_	<b>0.06</b> (0.04, 0.09)	<b>0.07</b> (0.05, 0.10)	_	<b>0.07</b> (0.05, 0.10)
Time Fixed Effects	Incl.	Incl.	Incl.	Incl.	Incl.	Incl.
City-Specific Time Fixed Effects	Incl.	Excl.	Excl.	Incl.	Excl.	Incl.
Community-Specific Time Fixed Effects	Incl.	Excl.	Excl.	Excl.	Incl.	Incl.
LMD	-950,255	-851,128	-850,403	-849,942	-849,485	-848,334

 Table 4. Estimation Results of Different Models

*Notes.* The 95% posterior confidence intervals for selected parameters are in parentheses. Estimates are bolded if the 95% posterior confidence intervals do not contain zero. For individual fixed effects, the mean and variance are computed post hoc based on individual draws. Excl., excluded; Incl., included.



Figure 4. Histograms of Median Posterior Estimates of Individual-Specific Social-Influence Parameters

median posterior estimates of the individual-specific social-influence parameters that are very similar to those of the full model 6. Overall, these patterns provide initial support that community-specific time fixed effects are useful for the identification of social influence, which we further explore in the next section.

We focus on the estimates of the full model (model 6), because it has the best model fit, and parameter estimates are qualitatively similar across models. First and foremost, we find a significant estimate of the susceptibility to social influence (mean: 0.91) that varies strongly across consumers (variance: 0.61),<sup>22</sup> as illustrated in Figure 4. Although the susceptibility to social-influence parameter is overall positive, some users are not affected by their friends' login decisions or have a slightly negative value. In particular, for 5.31% of all users, the median posterior estimates of the individual-specific social-influence parameters are negative. For the marketing variables, the regular promotion has a positive effect on logins (mean estimate: 0.09). Similarly, both targeted promotions significantly increase logins, and show comparable

effects (mean estimates: 0.06 and 0.07, respectively, for International Women's Day and the White Day).

Using these estimates, we quantified  $Influence_i$  that is, the social influence of user *i* on all other users in the network. To do so, we compared the current scenario with a counterfactual scenario where user *i* becomes inactive (i.e., the intrinsic utility is set to  $-\infty$ ). We did this for all time periods *T* in our data and subsequently computed the average across time periods, with the following equation:

$$Influence_{i} = \frac{1}{T} \sum_{t=1}^{T} \sum_{j \in N \setminus i} \left| \underbrace{\begin{array}{c} G_{\infty}^{j}(F(\Delta_{t})) \\ \text{Login probability of user } j \text{ at time } t, \\ \text{if user } i \text{ logs in as usual} \end{array} \right|^{- \underbrace{G_{\infty}^{j}(F(\Delta_{t}|V_{it} = -\infty))}_{\text{Login probability of user } j \text{ at time } t, \\ \text{if user } i \text{ becomes inactive}} \right|.$$
(11)

Using Equation (11), Figure 5 reports the results of the influence that each user has on other users in the network. As illustrated in Figure 5, logins of almost all users have a positive impact on login decisions of others, and only 51 users out of 25,418 users have negative influence (min influence: -0.40). This is consistent with our observations that users in the online game treated friendships very carefully for gameplay. Interestingly, the shape of Figure 5 is right skewed, with a small group of gamers having a disproportionally larger social influence (maximum social influence = 0.94 logins per quarter day). This shape is possibly driven by the distribution of the number of friends of users, which follows a power-law distribution as discussed in Section 4.1. In line with this, we find that social influence is positively correlated with degree centrality (r = 0.61 and p < 0.01).

#### 4.4.2. Validation of the Identification of Social Influence.

To identify social influence, we propose to include community-specific time fixed effects to proxy for exogenous "local shocks" that only affect users in the same community. Our estimation results support that community-specific time fixed effects indeed reduce the upward bias in the mean estimate of the socialinfluence parameters. However, the identification of the proposed structural model, as any structural model, relies on structural assumptions. To further validate the estimates of social-influence parameters of the structural model, we constructed an exogeneity test that utilizes instrumental variables (Hausman 1978, Rivers and Vuong 1988), which is executed as follows. With the estimation of the structural models in Table 4, we first computed the expected equilibrium login probabilities  $\hat{P}_t^*$  using Equation (5) and then derived the expected percentages of friends that are online (i.e.,  $\hat{A}\hat{P}_t^*$ ). We then used  $\hat{A}\hat{P}_t^*$  as the variable of social influence and estimated a reducedform model that is analogous to the linear-in-means specification (Manski 1993). The reduced-form model has the same specification as Equation (10), except that equilibrium login decisions are replaced by the predicted values (i.e.,  $\hat{A}\hat{P}_t^*$ ) from the structural estimations:

$$U(d_{it} = 1 | x_{it}, \hat{p}_{t}^{*}, \Theta_{i})$$

$$= \beta_{0i} + \lambda_{t} + \beta_{1i} x_{t} + \beta_{2} x_{it}^{\text{prom}} + \lambda_{t}^{\text{city}} x_{i}^{\text{city}} + \lambda_{t}^{\text{com}} x_{i}^{\text{com}}$$

$$+ \gamma_{i} \sum_{j \in N_{i}} \tilde{a}_{ij} \hat{p}_{jt}^{*} + \eta_{it}.$$
(12)

We assume that the error term  $\eta_{it}$  follows a standard normal distribution, and user *i* logs into the game at time period *t* (i.e.,  $d_{it} = 1$ ) if the latent utility is positive. Thus, the model corresponds to a probit model with the predicted percentages of friends that are online as the variable of social influence. Using this model, we examined whether the structural models are wellspecified with exogeneity tests on  $\tilde{A}\hat{P}_t^*$ .

Figure 5. The Distribution of Social Influence on Others Across Users



To construct the exogeneity test, we adopt a controlfunction approach by specifying a first-stage model with  $\tilde{A}\hat{P}_{t}^{*}$  as the dependent variable (Wooldridge 2010):

$$\sum_{j \in N_i} \tilde{a}_{ij} \hat{p}_{jt}^* = \alpha_{0i} + \lambda_t + \alpha_{1i} x_t + \alpha_2 \bar{x}_{it}^{\text{prom}} + \bar{\lambda}_t^{\text{city}} \bar{x}_i^{\text{city}} + \bar{\lambda}_t^{\text{com}} \bar{x}_i^{\text{com}} + \nu_{it}.$$
(13)

In Equation (13),  $\bar{x}_{it}^{\text{prom}} = \sum_{j \in N_i} \tilde{a}_{ij} x_{jt}^{\text{prom}}$  is the fraction of user *i*'s friends that are targeted during the targeted promotions. Similarly,  $\bar{x}_i^{\text{city}} = \sum_{j \in N_i} \tilde{a}_{ij} x_j^{\text{city}}$  and  $\bar{x}_i^{\text{com}} = \sum_{j \in N_i} \times$  $\tilde{a}_{ij} \mathbf{x}_i^{\text{com}}$  are the vectors containing the fraction of user *i*'s friends in each city and community, respectively. Following the control-function approach (Petrin and Train 2010), we assume the error terms of the first- and second-stage model to be correlated, with  $\eta_{it} = \theta v_{it} + e_{it}$ . The test for the exogeneity of  $\widehat{AP}_t^*$  is equivalent to the test whether  $\theta = 0$  (Wooldridge 2010). The main challenge of implementing the test is to identify an instrumental variable that is excluded from user i's login decisions, but indirectly affects logins of the user through his friends. In our data set, we observed two targeted promotions and used  $\bar{x}_{it}^{\text{prom}}$ , the fraction of friends that are targeted by the two promotions, as instrumental variables. Because all users were informed about targeted promotions beforehand, rational users would integrate these promotions into their expectations about the login decisions of their friends, making targeted promotions a valuable instrument.

Next, we discuss how the proposed instruments,  $\bar{x}_{_{it}}^{\mathrm{prom}}$ , satisfy three conditions, as discussed by Imbens (2014): (1) exogeneity, (2) exclusion restriction, and (3) relevance. First, the proposed instruments are exogenous, as the two targeted promotions depend on the timing and nature of the public events, which are exogenous. Second, the proposed instruments are excluded from the login decisions of user *i*. This is because we control for targeted promotions  $x_{it}^{prom}$  of user *i* and time fixed effects  $\lambda_t$  as explanatory variables in Equation (12). Any remnant impact of targeted promotions on the focal user is through his friends. Third, to examine the relevance of the instruments, we ran two first-stage regressions, as in Equation (13), with the endogenous variables (i.e.,  $\hat{A}\hat{P}_{t}^{*}$ ) calculated from the full model (model 6 in Table 4) and one of the two instruments. As seen from models (b) and (c) in Table 5, the parameter estimates of the instruments are positive and significant (estimates: 0.009 and 0.008, respectively, for the targeted promotion on International Women's Day and White Day, with none of the posterior draws negative). To further examine whether the instruments are relevant, we ran a baseline model (model a) without instruments. The differences between LMDs (i.e., logged Bayes factors) of both regressions and the baseline model are 87.21 and 57.25, respectively, which are very strong evidence that the instruments are relevant (Kass and Raftery 1995). In addition, we also ran likelihood-ratio

	Model (a)	Model (b)	Model (c)	Model (d)
Variables	No instrument	With Women's Day instrument	With White Day instrument	With both instruments
Instruments (Fraction of Friends Targeted)				
Promotion (Women's Day)	_	<b>0.009</b> (0.004, 0.014)	_	0.009 (0.005, 0.013)
Promotion (White Day)	_	_	0.008 (0.001, 0.016)	0.008 (0.002, 0.014)
Average Individual Fixed Effects of Friends				
Post hoc Mean	0.52	0.52	0.52	0.53
Post hoc Variance	0.45	0.45	0.45	0.45
Average Time Effects of Friends				
Regular Promotion	0.02 (0.02, 0.02)	<b>0.02</b> (0.02, 0.03)	0.02 (0.01, 0.02)	<b>0.02</b> (0.00, 0.02)
Time Fixed Effects	Incl.	Incl.	Incl.	Incl.
Average City-specific Time Fixed Effects of Friends	Incl.	Incl.	Incl.	Incl.
Average Community-specific Time Fixed Effects of Friends	Incl.	Incl.	Incl.	Incl.
LMD	-762,612	-762,524	-762,554	-762,436
Bayes Factor (log)	—	87.21	57.25	145.38
F statistics	—	351.19	252.82	493.86

**Table 5.** Estimation Results of the First-Stage Model

*Notes.* The dependent variable is the average expected logins of friends as computed based on the structural model (full model 6 in Table 4). Models (a)–(c) are estimated separately from the second-stage model, and model (d) is estimated jointly with the second-stage model. The 95% posterior confidence intervals for selected parameters are in parentheses. Estimates are bolded if the 95% posterior confidence intervals do not contain zero. *F* statistics and Bayes factors are calculated with model (a) as the baseline model. For individual fixed effects, the mean and variance are computed post hoc based on individual draws. Excl., excluded; Incl., included.

For the control-function approach, we jointly estimated Equations (12) and (13). Table 6 presents the results from the control-function estimation, as well as the direct estimation. We tested specifications of structural models 4–6 in Table 4,<sup>23</sup> to examine the usefulness of communities in the structural estimation. The estimation results reveal the following. First, the full model 6 shows no evidence of endogeneity, as the estimate of covariance  $\theta$  is insignificant, with the posterior 95% confidence interval containing zero. Second, the estimate of covariance  $\theta$  based on the variable of social influence from model 5 is also insignificant (95% posterior confidence interval: [-0.38, (0.06]), whereas that from model 4 is significant (95%) posterior confidence interval: [-0.48, -0.38]). These results support the specification of model 5, which only includes community-specific time fixed effects as excluded variables. In contrast, model 4, which includes targeted promotions and city-specific time fixed effects, is not supported by the exogeneity test. Overall, these results further support the use of communities to construct excluded variables to identify social influence.

#### 4.5. Implications for Targeting

An advantage of the structural model is that it allows researchers to perform counterfactual simulations to evaluate the effectiveness of new strategies not covered in the data. In this section, we illustrate the potential managerial usefulness of the proposed structural model of social influence. As discussed in Section 4.1, keeping users active is vital for the survival of online games, as users' engagement generates revenue through in-game purchases. To increase the activity level of users, the online game could leverage the "social multiplier effect" by targeting a small group of influential users. Traditionally, to increase engagement, companies mainly target consumers based on their usage patterns (Ballings and Van den Poel 2015) or strategic positions in the network, such as degree centrality (Hinz et al. 2011, Chen et al. 2017). As a novel enhancement to these approaches, we propose to leverage how users respond to social influence across time and illustrate how the online game can improve the effectiveness of targeting. The important role of timing in marketing communications has been recognized by both marketing practitioners

and academics (Drèze and Bonfrer 2009). We are all used to receive telemarketing calls right before dinnertime and an inflow of promotional emails to our mailboxes between 9 a.m. and 10 a.m. in the morning. Clearly, these times are not chosen randomly, but reflect marketers' anticipation of returns on their efforts. Likewise, in the domain of social networks, the timing of reaching out to influential consumers may play a critical role.

We argue that when choosing the optimal timing in network settings, marketers need to consider several factors. First, marketers need to consider the probability that the marketing intervention changes the target's behavior. This is a well-known principle of marketing communication, as it would be wasteful to stimulate an action of consumers who are highly likely to take the action without additional stimulation, or are highly unlikely to react to the stimulation. Second, depending on the time of the intervention, a consumer's response (i.e., changes in behaviors) may vary. In the online game, there are times of targeting that are preferable to some users, but less popular among others. Finally, how users respond to their friends may also depend on the time of the day. Hence, to optimally select the right targets with the highest potential impact on the network, it is important to have a model that captures all of the above dimensions. The model developed in this paper may serve this purpose, as it allows one to predict the individual-specific responses to a certain amount of direct and/or social stimulation that is received at a given time.

The intuition behind the proposed time-/socialawareness targeting approach is as follows. At a given point in time, the most promising targets are the network members who (1) are likely to respond (i.e., to change their behavior) to the *direct* marketing stimulation, and (2) have the most direct (i.e., friends) and indirect connections (i.e., friends of friends) who are likely to respond to the *social* stimulation that is resulted from the target's change of actions. We explore the performance contribution of these factors through policy simulations. To obtain the performance of a targeting decision, we assume that the company is able to stimulate the intrinsic login utility of a group of users  $S \subset N \equiv \{1, 2, ..., N\}$  (i.e., *S* is a subset of all network users, and *N* is the set of all users). For each user s in set S, we assume that the intrinsic utility increases by  $\delta_s$  due to the targeting efforts. We set  $\delta_s$ equal to the estimated individual-specific response of user *s* to the regular promotion from model 6. Following the intuition in Equation (11), the performance at time *t* of targeting users in subset *S* with marketing stimulation  $\delta^{S}$  (i.e.,  $\Pi_{S_{i}}^{(\delta^{S})}$ ) can be calculated as the

	Model 4 without community		Model 5 without city- targeted promotions		Model 6 full model	
Variables	CF-Probit	Probit	CF-Probit	Probit	CF-Probit	Probit
Social Influence						
Mean	<b>1.03</b> (1.01, 1.05)	<b>1.11</b> (1.08, 1.13)	<b>0.92</b> (0.91, 0.94)	<b>0.92</b> (0.91, 0.94)	<b>0.91</b> (0.86, 0.96)	<b>0.91</b> (0.89, 0.94)
Variance	<b>0.78</b> (0.76, 0.81)	<b>0.80</b> (0.78, 0.83)	<b>0.67</b> (0.65, 0.69)	<b>0.67</b> (0.65, 0.69)	<b>0.66</b> (0.64, 0.68)	<b>0.66</b> (0.64, 0.67)
Individual Fixed Effects						
Post hoc mean	-1.63	-1.66	-1.52	-1.51	-1.52	-1.52
Post hoc variance	0.34	0.33	0.33	0.33	0.38	0.38
Promotions						
Regular Promotion	<b>0.06</b> (0.01, 0.12)	<b>0.06</b> (0.01, 0.11)	<b>0.06</b> (0.01, 0.11)	<b>0.06</b> (0.01, 0.11)	<b>0.08</b> (0.03, 0.14)	<b>0.08</b> (0.01, 0.15)
Promotion (Women's Day)	<b>0.06</b> (0.04, 0.08)	<b>0.06</b> (0.04, 0.08)	<b>0.06</b> (0.04, 0.09)	<b>0.06</b> (0.04, 0.09)	<b>0.06</b> (0.04, 0.08)	<b>0.06</b> (0.04, 0.08)
Promotion (White Day)	<b>0.07</b> (0.05, 0.10)	<b>0.07</b> (0.05, 0.10)	<b>0.08</b> (0.05, 0.10)	<b>0.08</b> (0.05, 0.10)	<b>0.08</b> (0.05, 0.10)	<b>0.08</b> (0.05, 0.10)
Time Fixed Effects	Incl.	Incl.	Incl.	Incl.	Incl.	Incl.
City-Specific Time Fixed Effects	Incl.	Incl.	Incl.	Incl.	Incl.	Incl.
Community-Specific Time Fixed Effects	Incl.	Incl.	Incl.	Incl.	Incl.	Incl.
Covariance $(\theta)$	<b>-0.43</b> (-0.48, -0.38)	—	-0.16 (-0.38, 0.06)	—	-0.10 (-0.31, 0.10)	—

#### Table 6. Testing the Specifications of Structural Models

*Notes.* "CF-Probit" stands for the second-stage probit models that are estimated with the control-function approach. The 95% posterior confidence intervals for selected parameters are in parentheses. Estimates are bolded if the 95% posterior confidence intervals do not contain zero. For individual fixed effects, the mean and variance are computed post hoc based on individual draws. Excl., excluded; Incl., included.

overall increase of logins of the network, compared with the situation without targeting:



where  $\delta^{S}$  is a ( $N \times 1$ )-vector indicating the change in intrinsic utilities due to targeting (i.e., element *i* of  $\delta^{S}$  equals  $\delta_{i}$  if  $i \in S$ , and zero otherwise), and  $\Delta_{t}$  is a vector of intrinsic utilities without targeting.

**4.5.1. The Decision of "Whom to Target.**" To examine the decision of "whom to target", we compared the following four targeting approaches: (1) "Homogeneous Parameters," (2) "Hub," (3) "Responder," and (4) "Influencer." We compared these heuristics using Equation (14) during the time of the regular

promotion (quarter days 85–104 or days 22–26). The Homogeneous Parameters approach targets users that maximize Equation (14), but uses the estimates of model 1 (Homogeneous Parameters; see Table 4). The Hub approach focuses on the dense regions of the social network and targets those users with the most connections. The Responder approach targets users who are most likely to change their logins if incentivized (i.e., with the highest increase in login probabilities if targeted with the promotion). The Influencer approach focuses on both responsiveness and network positions and targets those users that maximize Equation (14) using the parameter estimates of the full model 6. In all four targeting approaches, we selected 1,000 users.

The targeting results of the four approaches are presented in Table 7. The inclusion of individualspecific parameters strongly improves the targeting performance, as the expected performance of the Homogeneous Parameters approach is the worst. This highlights the importance of extending the traditional model (e.g., Lee et al. 2014) with individualspecific parameters. Interestingly, targeting with the Hub approach leads to much worse results than the Responder approach (360.72 versus 3,627.50), which highlights the importance of considering the responses of users to the targeting incentives. In particular, the

Targeting approaches	Targeting performances	Improvement over homogeneous approach
Homogeneous parameters	360.72 (18.04)	_
Hub	562.38 (28.12)	55.91%
Responder	3,627.50 (181.38)	905.64%
Influencer	4,268.87 (213.44)	1,083.44%

 Table 7. Comparing Different Targeting Approaches: Who to Target?

*Notes.* Targeting performances indicate the expected total number of additional logins during quarter days 85–104 by targeting 1,000 users. The numbers in parentheses are the corresponding average expected additional logins per quarter day. The Responder approach optimally selects targets based on their responsiveness. The Hub approach optimally selects targets based on degree centrality. The Influencer approach optimally selects users based on their responsiveness and the responsiveness of connected users in their network. The Homogeneous Parameters approach is similar to the Influencer approach, but uses the model estimates with homogeneous parameters (model 1 in Table 4).

1,000 most connected users have significantly lower average responsiveness than the least connected 1,000 users (mean: 0.04 versus 0.08, with *T*-value: -10.12 and p < 0.00). This result is in line with previous research (e.g., Gelper et al. 2020) and illustrates the potential risks of the commonly used Hub approach, as highly connected users may be nonresponsive to marketing efforts. More importantly, the Influencer approach, which considers both direct responsive-ness and connectedness of users, significantly outperforms all other approaches (expected targeting performance = 4,268.87 versus 3,627.50, respectively, for the Influencer approaches).

**4.5.2. The Decision of "When to Target.**" The timing of marketing interventions plays an important role in targeting, as consumers' responsiveness is expected to vary across time. In an environment where users are connected by a social network, an additional level of complexity arises from potential social responsiveness that may also be time-dependent. To explore the effect of timing on network targeting, we assume that the gaming company faces a decision of choosing one of four quarters of the day (12 a.m.–6 a.m., 6 a.m.–12 p.m., 12 p.m.–6 p.m., and 6 p.m.–12 a.m.), during the five consecutive days of the regular

promotion. Based on this decision problem, we compared two scenarios with the Influencer approach that does not consider timing.<sup>24</sup> In the first scenario, named "Uniform Timing," the company selects an optimal quarter to stimulate all of the targeted users. In the second scenario, named "Personalized Timing," each targeted user is stimulated at a personalized "optimal" quarter, which depends on the responsiveness of the user, as well as the responsiveness of the user's friends.<sup>25</sup>

The simulation results are reported in Table 8. Compared with the Influencer approach that does not consider timing, the Uniform Timing approach improves the targeting performance by 35.06% (targeting performances: 1,441.38 and 1,067.22, respectively, for the Uniform Timing and Influencer approaches). Moreover, using a personalized promotion schedule further improves the targeting performance by 15.61% (targeting performance of the Personalized Timing: 1,607.96).

Finally, in Table 9, we compare the optimal sets of users that are targeted at each quarter with a Uniform Timing approach. As illustrated in this table, the optimal selection of targeted users differs significantly across the quarters of the day, with overlaps ranging from 24.81% (quarters 1 and 3) to 46.95%

Table 8. Comparing Different Targeting Approaches: When to Target?

Targeting approaches	Targeting performances	Improvement over influencer approach
Influencer (no timing)	1,067.22 (213.44)	_
Uniform timing	1,441.38 (288.28)	35.06%
Personalized timing	1,607.96 (321.59)	50.67%

*Notes.* Targeting performances indicate the total number of additional logins by targeting 1,000 users on five quarter-days, with each quarter-day on one of the five consecutive days. The numbers in parentheses are the corresponding average expected additional logins per quarter day. Both the Uniform Timing and Personalized Timing approaches target the same set of users as the Influencer approach. The Uniform Timing targets all users in the overall optimal quarter (i.e., quarter 4), while the Influencer approach randomly selects a time period. The Personalized Timing approach personalizes the quarter of the day for each individual.

Quarters of the day	Quarter 1	Quarter 2	Quarter 3	Quarter 4
Quarter 1 (00:00–06:00)	100%	34.79%	24.81%	25.90%
Quarter 2 (06:00–12:00)		100%	34.04%	32.45%
Quarter 3 (12:00–18:00)			100%	46.95%
Quarter 4 (18:00-24:00)				100%

Table 9. The Overlap of Selected Users Across Quarters

*Notes.* Overlaps indicate the percentage of users that are present in both optimal target sets  $S_q$  and  $S_{q,r}$  with q and q' corresponding to two quarters that are targeted with the Uniform Timing approach. For instance, if q = 1 and q = 2, the overlap equals 34.79%.

(quarters 3 and 4). In other words, the set of early bird influencers are quite different from the night owls. To the best of our knowledge, this important targeting factor has received limited attention in the extant marketing literature on targeting under social influence.

#### 5. Discussion

The discrete choice model of social influence has drawn much attention from both marketing academics and practitioners, because of the increasing need to understand how consumers influence the decisions of one another. As such, researchers demand game-theoretic choice models that allow for individual-specific parameters to capture the heterogenous responses of consumers to social-influence and marketing activities. Applying such models to large social networks is challenging due to high computational demands and the difficulties of obtaining excluded variables for identification. In this paper, we provided novel solutions to address these two challenges. First, following Imai et al. (2009), we proposed a stochastic Bayesian estimation procedure that significantly reduced the computational demands. Simulation studies showed that computational complexity reduced from polynomial to linear in the network size, while all parameters were effectively recovered. Second, as a novel identification strategy, we recommended constructing communityspecific time fixed effects as excluded variables based on community-detection algorithms. The proposed identification strategy builds on the idea that people in different communities are likely exposed to different external shocks. The empirical analysis of login decisions in an online game validated this identification strategy with exogeneity tests that utilized targeted promotions as instrumental variables.

With an empirical application that involves login decisions of 25,000 users in an online game, we further demonstrated the managerial usefulness of the proposed methodology. First, although it is nearly impossible for traditional estimation methods to apply the proposed structural model to such a large network, the stochastic estimation procedure completed the estimation in reasonable time. Second, based on the model, we proposed a procedure for marketers to quantify the social influence of users with counterfactual simulations, which illustrated the importance of incorporating individual-specific parameters. Finally, counterfactual simulations showed that the model can significantly improve targeting decisions compared with traditional approaches that focus on network positions or responsiveness to marketing stimuli. As a novel insight, these analyses highlighted the importance of timing in targeting individual consumers to leverage social influence, as the optimal timing varies substantially across consumers.

Future research can apply our model to other empirical scenarios, where social influence is a significant driver of consumer choices. For instance, Online Appendix I applies the proposed methodology to the login decisions of users in an online social network. In this empirical application, we did not observe any excluded variables and further proved the usefulness of the proposed identification strategy based on latent communities. Interestingly, in this application, we find a small group of users with negative responses to social influence, which may be explained by psychological reactance to social influence. For instance, recent reports suggested that some users avoid their online friends, so that they are not disturbed by chat requests.<sup>26</sup> In addition, the setting of donations is of potential interest, given that donors' decisions are often driven by social pressure that is originated from the expectations about the donation decisions of peers. Researchers can also study the purchase decision of luxury products, where a purchase becomes less attractive if consumers expect others to make the same purchase. Finally, the estimation method can also be applied to interactions between other market agents, such as the competition between firms (e.g., Seim 2006 and Vitorino 2012).

There are several opportunities to extend the proposed model. First, the model can be extended to multiple discrete-choice scenarios, where the equilibrium is defined by a matrix of choice probabilities, with each column corresponding to one of the options. Second, future research can incorporate network formation. Although in the current version, we include community-specific shocks, a network formation model may be useful to control for the endogeneity of social networks.<sup>27</sup> Moreover, a dynamic network formation model can extend our model to social networks that are in the growing or declining stage of the lifecycle. Third, integrating a community-detection procedure into the structural model to guide the selection of the proper communities is another fruitful direction for future research. Given the efficiency of the stochastic estimation method, such a practice may be feasible for very large networks. Moreover, considering behavioral similarities may improve the accuracy of existing community-detection algorithms that mainly rely on network structure (Fortunato 2010), sometimes complemented with user characteristics (Yang et al. 2013). Finally, in some empirical settings, it is conceivable that users are forwardlooking, such that they consider their friends' future reactions to their current decisions. The inclusion of a forward-looking component into the model changes the static network game to a dynamic one. To apply the stochastic estimation procedure, it is critical to derive the conditions under which the dynamic network game admits a unique solution.

In sum, we illustrated how the stochastic Bayesian estimation procedure, in combination with community detection, allows researchers to efficiently and effectively estimate social influence in large social networks. The proposed methodology is flexible and does not require researchers to collect valid excluded variables that are often difficult to obtain in empirical applications. We hope that the current paper inspires researchers to apply structural choice models of social influence to large social networks.

#### Acknowledgments

The authors thank the editorial team for valuable suggestions throughout the review process; Cheng Zhang (Fudan University) and Xueming Luo (Temple University) for sharing the main data in the paper; and Andrew Ching (Johns Hopkins University) for valuable suggestions on an earlier version of this paper.

#### Endnotes

<sup>1</sup> Because we allow the susceptibility to social-influence parameter ( $\gamma_i$ ) to be individual-specific, the average and aggregate models are mathematically equivalent. To illustrate this, define  $\tilde{a}_{ij} = a_{ij}/N_i$ , the row-normalized connection between individuals *i* and *j*, with  $deg_i$  the (weighted) degree of individual *i* or summation of row *i* of matrix *A*. Using the right-hand side of Equation (2), we have  $\gamma_i \sum_{j \in N_i} a_{ij} \mu_{ijt}^e (d_{jt} = 1) = \gamma_i deg_i \sum_{j \in N_i} \tilde{a}_{ij} \mu_{ijt}^e (d_{jt} = 1)$ . For an aggregate model  $(a_{ij})$  and its average counterpart  $(\tilde{a}_{ij})$ , we must have  $\gamma_i^{\text{aggregate}} \times deg_i = \gamma_i^{\text{average}}$ . However, in empirical applications, it is common to assume a distribution for  $\gamma_i$ . Hence, in such scenarios, the aggregate and average model differ in the distributional assumptions of  $\gamma_i$ . In our empirical applications, we used model fit statistics to determine which specification is better.

<sup>2</sup> Under this assumption, the model can still be applied to social networks that vary across time if consumers are not strategically

forming and/or dissolving connections in consideration of the choice decisions.

<sup>3</sup> The rational expectation assumption implies that, to form a correct belief about a friend *j*, consumer *i* needs to form correct expectations about the choice probabilities of *j*'s friends. This may be a strong assumption. Therefore, we follow the theory of learning in games (Fudenberg and Levine 1998) and assume that consumers have already learned how to form expectations, and the game always reaches the locally stable equilibrium (see Proposition 1). This is also in line with the assumption that consumers are myopic and that the data observed in each time period are from an independent realization of the game.

<sup>4</sup> The condition for a unique equilibrium in the case of homogenous parameters is  $|\gamma| < 1/(max(deg_i) \cdot f^{max})$ . Note that in both the homogenous and individual-specific parameter cases, these conditions are sufficient, but not necessary.

<sup>5</sup> The inclusion of individual fixed effects may lead to a finite sample bias when the number of time periods *T* is small. However, as shown in Greene (2004), the bias decreases when *T* increases, and recent research has shown that the bias becomes marginal when *T* is large (see, for example, Ibanez et al. 2018 and Stafford 2015). Our empirical data have 120 time periods, which allow for the inclusion of individual fixed effects. To examine whether *T* of our empirical data are sufficient, we followed Greene (2004) and performed simulation studies with data sets that have 100 time periods. The simulation results in Section 3.4 demonstrate that individual fixed effects, as well as other parameters, are well-recovered.

<sup>6</sup> In our empirical application, we find support for this conjecture, as we were able to compare excluded variables based on community detection with those constructed from observed geo-locations.

<sup>7</sup> Some previous research assumed that only one equilibrium is realized in the data and selected one equilibrium out of all possible equilibria with criteria such as Pareto optimality (e.g., Hartmann 2010).

<sup>8</sup>Notice that with individual-specific parameters considered, the algorithm requires large memory space to store previous draws. In our application, instead of focusing on parameter vectors  $\Theta$ , we use draws of latent utilities ( $z_i$ ; see Online Appendix B.2). Such a strategy reduces the array size from  $R \times N \times (K + 1)$  to  $R \times N \times 1$  and, thus, significantly reduces memory requirements.

<sup>9</sup> We further prove in Online Appendix B.1 that under assumptions of the proposed model, the parameter draws from the proposed estimator converge to the true posterior distributions. Specifically, we show that the proposed estimator meets the conditions to apply theorem 2 of Imai et al. (2009).

<sup>10</sup> In the estimation, we assume that the community structure is known. As illustrated by Girvan and Newman (2002) and Newman (2006), the spectral bisection community-detection method accurately recovered the community structure in networks that were generated by using this simulation approach. In our empirical applications, we find that the estimation results are robust for different community structures detected by different methods, including the spectral bisection method (see Online Appendices E and I.5.2).

<sup>11</sup>The estimates of the other parameters are almost identical for both estimation methods. Because our focus is on social influence, we only report the estimates of social-influence parameters in Table 1. For the estimation of individual fixed effects, time fixed effects, and community-specific time fixed effects, please refer to Online Appendix C.

<sup>12</sup>We also estimated models that do not control for individual fixed effects, which, as expected, also result in upward biased estimates of social influence (see Online Appendix C.3).

<sup>13</sup>https://www.statista.com/study/39310/video-games-2018/ (accessed June 2, 2020).

7597

<sup>14</sup>We acknowledge that the analysis does not fully rule out the possibility that users are forward-looking. For instance, if users are satiated, they may decide to postpone logging into the game, if they expect that their friends are more likely to login in the future (we thank an anonymous reviewer for pointing out this possible source of forward-looking tendencies). Therefore, extending the model to incorporate forward looking is an important direction for future research, which we discuss in the final section.

<sup>15</sup>In the community-detection analysis, we compared the communities across daily networks and did not observe any differences. In our estimation and subsequent analyses, we always take into consideration the daily networks.

<sup>16</sup> White Day is a popular event in many East Asian countries and is celebrated one month after Valentine's Day. During this event, men are expected to give gifts to their loved ones. The promotion incentivized male users to give gifts to female users in return for valuable in-game items. See Online Appendix D for more details of the promotions.

<sup>17</sup>We also compared our estimation results with communities obtained using different algorithms and found that the estimation results are robust across different algorithms (see Online Appendix E).

<sup>18</sup>Collinearity diagnostics did not reveal multicollinearity problems, with the maximum variance inflation factor equal to 3.84. Also, in cases when a large number of communities are detected, regularization techniques (e.g., LASSO) could be used to address the potential overfitting problem caused by the high dimensionality in the predictors (e.g., see Gillen et al. 2019 and Kang et al. 2016).

<sup>19</sup>To check the plausibility of the independence of error terms  $\varepsilon_{it}$ , we tested the residuals of this model (full model 6 discussed below) for each community. We used Breusch–Godfrey test and found no evidence of serial correlation over time, with the *p*-values of the 125 communities ranging from 0.15 and 0.83. In addition, we used Moran's *I* to test for the possible spatial correlation between users within each community and found no evidence of spatial correlation, with *p*-values ranging from 0.22 to 0.99.

<sup>20</sup> To examine the validity of these restrictions on the susceptibility to social-influence parameters, we also estimated the model without the restrictions. Our estimation results are robust, and the restricted model fits the data better. The results further support the validity of our model specification (see Online Appendix F for more details). In addition, we estimated a model with the aggregate specification—that is, using the original adjacency matrix *A* instead of the row-normalized matrix  $\tilde{A}$ . The LMD of this model is -848,472, smaller than that of the average model (*LMD* = -848,334 for model 6 of Table 4). These results support the specification of the average model.

<sup>21</sup> To include the regular promotion, we must set another time fixed effect to zero. Because the regular promotion lasted from quarter-day 85 to 104, we set quarter-day 85 to zero.

<sup>22</sup>We found a weak positive correlation between susceptibility to social influence and degree centrality (r = 0.23, p < 0.01). Therefore, being susceptible to social influence is not necessarily driven by the number of friends, as some highly susceptible users may have only a few friends.

<sup>23</sup> For the first-stage estimation results of model 4 and model 5, please see Online Appendix H.

<sup>24</sup>The Influencer approach effectively assumes that the company randomly selected one of the four quarters to target users, resulting in 1/4 of the targeting performance in Table 7 as the expected performance of this approach.

<sup>25</sup>Note that the selection of targets in both approaches follows the Influencer strategy, but are enhanced with the time dimension. The difference between the two scenarios is that in the Personalized Timing case, the restriction of setting the same timing of the marketing intervention for all users is relaxed.

#### <sup>26</sup> https://www.businessinsider.com/how-to-appear-offline-on-facebook (accessed April 4, 2020).

<sup>27</sup> There is limited research integrating network formation into network games. Goldsmith-Pinkham and Imbens (2013) developed an approach that integrates a logistic regression model for social-network formation into a linear-in-means model. However, it is challenging to scale this model to large social networks, given that the total number of potential links between consumers grows quadratically. With 25,418 users in our online gaming data, the logistic regression needs to be applied to more than 646 million observations.

#### References

- Ahn D-Y, Duan JA, Mela CF (2015) Managing user-generated content: A dynamic rational expectations equilibrium approach. *Marketing Sci.* 35(2):284–303.
- Ascarza E, Ebbes P, Netzer O, Danielson M (2017) Beyond the target customer: Social effects of customer relationship management campaigns. J. Marketing Res. 54(3):347–363.
- Bainbridge WS (2007) The scientific research potential of virtual worlds. *Science* 317(5837):472–476.
- Bajari P, Hong H, Krainer J, Nekipelov D (2010) Estimating static models of strategic interactions. J. Bus. Econom. Statist. 28(4): 469–482.
- Ballings M, Van den Poel D (2015) CRM in social media: Predicting increases in Facebook usage frequency. *Eur. J. Oper. Res.* 244(1):248–260.
- Blondel VD, Guillaume J-L, Lambiotte R, Lefebvre E (2008) Fast unfolding of communities in large networks. J. Statist. Mechanics Theory Exp. (10):P10008.
- Brock WA, Durlauf SN (2001) Discrete choice with social interactions. *Rev. Econom. Stud.* 68(2):235–260.
- Carlsson H, Damme EV (1993) Global games and equilibrium selection. *Econometrica* 61(5):989–1018.
- Chen X, van der Lans R, Phan TQ (2017) Uncovering the importance of relationship characteristics in social networks: Implications for seeding strategies. J. Marketing Res. 54(2):187–201.
- Csardi G, Nepusz T (2006) The igraph software package for complex network research. *Internat. J. Complex Systems* 1695(5):1–9.
- De Giorgi G, Pellizzari M, Redaelli S (2010) Identification of social interactions through partially overlapping peer groups. *Amer. Econom. J. Appl. Econom.* 2(2):241–275.
- Drèze X, Bonfrer A (2009) Moving from customer lifetime value to customer equity. *Quant. Marketing Econom.* 7(3):289–320.
- Ellickson PB, Misra S (2011) Structural workshop paper-estimating discrete games. *Marketing Sci.* 30(6):997–1010.
- Fortunato S (2010) Community detection in graphs. *Phys. Rep.* 3(486):75–174.
- Fudenberg D, Levine DK (1998) *The Theory of Learning in Games* (MIT Press Books, Cambridge, MA).
- Galeotti A, Goyal S, Jackson MO, Vega-Redondo F, Yariv L (2010) Network games. *Rev. Econom. Stud.* 77(1):218–244.
- Gelper S, van der Lans R, van Bruggen G (2020) Competition for attention in online social networks: Implications for seeding strategies. *Management Sci.*, ePub ahead of print June 30, https:// doi.org/10.1287/mnsc.2019.3564.
- Gillen BJ, Montero S, Moon HR, Shum M (2019) BLP-2LASSO for aggregate discrete choice models with rich covariates. *Econom. J.* 22(3):262–281.
- Girvan M, Newman ME (2002) Community structure in social and biological networks. Proc. National Acad. Sci. USA 99(12):7821–7826.
- Glaeser E, Scheinkman J (2001) Measuring social interactions. Durlauf SN, Young PH, eds. *Social Dynamics* (MIT Press, Cambridge, MA), 83–132.
- Goldsmith-Pinkham P, Imbens GW (2013) Social networks and the identification of peer effects. J. Bus. Econom. Statist. 31(3):253–264.

- Greene W (2004) The behaviour of the maximum likelihood estimator of limited dependent variable models in the presence of fixed effects. *Econom. J.* 7(1):98–119.
- Gruber J, Köszegi B (2001) Is addiction "rational"? Theory and evidence. *Quart. J. Econom.* 116(4):1261–1303.
- Gupta S, Mela CF (2008) What is a free customer worth? Armchair calculations of nonpaying customers' value can lead to flawed strategies. *Harvard Bus. Rev.* 86(11):102–109.
- Gupta S, Hanssens D, Hardie B, Kahn W, Kumar V, Lin N, Ravishanker N, Sriram S (2006) Modeling customer lifetime value. J. Service Res. 9(2):139–155.
- Hansen LP, Sargent TJ, Heaton J, Marcet A, Roberds W (1991) *Rational Expectations Econometrics* (Westview Press, Boulder, CO).
- Hartmann W (2010) Demand estimation with social interactions and the implications for targeted marketing. *Marketing Sci.* 29(4): 585–601.
- Hausman J (1978) Specification tests in econometrics. *Econometrica* 46(6):1251–1272.
- Hinz O, Skiera B, Barrot C, Becker JU (2011) Seeding strategies for viral marketing: An empirical comparison. *J. Marketing* 75(6): 55–71.
- Hogan JE, Lemon KN, Libai B (2003) Quantifying the ripple: Word-ofmouth and advertising effectiveness. J. Advertising Res. 44(3): 271–280.
- Hotz VJ, Miller RA (1993) Conditional choice probabilities and the estimation of dynamic models. *Rev. Econom. Stud.* 60(3):497–529.
- Ibanez MR, Clark JR, Huckman RS, Staats BR (2018) Discretionary task ordering: Queue management in radiological services. *Management Sci.* 64(9):4389–4407.
- Imai S, Jain N, Ching A (2009) Bayesian estimation of dynamic discrete choice models. *Econometrica* 77(6):1865–1899.
- Imbens GW (2014) Instrumental variables: An econometrician's perspective. *Statist. Sci.* 29(3):323–358.
- Kang H, Zhang A, Cai TT, Small DS (2016) Instrumental variables estimation with some invalid instruments and its application to Mendelian randomization. J. Amer. Statist. Assoc. 111(513): 132–144.
- Kass RE, Raftery AE (1995) Bayes factors. J. Amer. Statist. Assoc. 90(430):773–795.
- Kumar V, Reinartz W (2012) Customer Relationship Management: Concept, Strategy, and Tools (Springer Science & Business Media, Berlin).
- Lancichinetti A, Fortunato S (2009) Community detection algorithms: A comparative analysis. *Phys. Rev. E* 80(5):056117.
- Lee L-f, Li J, Lin X (2014) Binary choice models with social network under heterogeneous rational expectations. *Rev. Econom. Statist.* 96(3):402–417.
- Leskovec J, Lang KJ, Mahoney M (2010) Empirical comparison of algorithms for network community detection *Proc. 19th Internat. Conf. World Wide Web* (Association for Computing Machinery, New York), 631–640.
- Lin X (2010) Identifying peer effects in student academic achievement by spatial autoregressive models with group unobservables. *J. Labor Econom.* 28(4):825–860.
- Magnac T, Thesmar D (2002) Identifying dynamic discrete decision processes. *Econometrica* 70(2):801–816.
- Manski CF (1993) Identification of endogenous social effects: The reflection problem. *Rev. Econom. Stud.* 60(3):531–544.
- Mazzeo MJ (2002) Product choice and oligopoly market structure. *RAND J. Econom.* 33(2):221–242.
- Misra S (2013) Markov chain Monte Carlo for incomplete information discrete games. *Quant. Marketing Econom.* 11(1):117–153.

- Myerson RB (1978) Refinements of the Nash equilibrium concept. Internat. J. Game Theory 7(2):73–80.
- Narayanan S, Nair HS (2013) Estimating causal installed-base effects: A bias-correction approach. J. Marketing Res. 50(1):70–94.
- Newman M (2006) Modularity and community structure in networks. *Proc. Natl. Acad. Sci. USA* 103(23):8577–8582.
- Newman ME, Clauset A (2016) Structure and inference in annotated networks. *Nat. Commun.* 7(1):1–11.
- Nicoletti C, Salvanes KG, Tominey E (2018) The family peer effect on mothers' labor supply. *Amer. Econom. J. Appl. Econom.* 10(3): 206–234.
- Petrin AK, Train K (2010) A control function approach to endogeneity in consumer choice models. J. Marketing Res. 47(1):3–13.
- Rivers D, Vuong Q (1988) Limited information estimators and exogeneity tests for simultaneous probit models. *J. Econometrics* 39(3):347–366.
- Ronhovde P, Nussinov Z (2009) Multiresolution community detection for megascale networks by information-based replica correlations. *Phys. Rev. E* 80(1):016109.
- Rosvall M, Bergstrom CT (2008) Maps of random walks on complex networks reveal community structure. *Proc. Natl. Acad. Sci. USA* 105(4):1118–1123.
- Rust JP (1987) Optimal replacement of GMC bus engines: An empirical model of Harold Zurcher. *Econometrica* 55(5):999–1033.
- Scott DW, Sain SR (2005) Multi-dimensional density estimation. Rao CR, Wegman EJ, Solka JL, eds. *Handbook of Statistics*, vol. 23 (North-Holland Publishing Co., Amsterdam), 229–261.
- Seim K (2006) An empirical model of firm entry with endogenous product-type choices. *RAND J. Econom.* 37(3):619–640.
- Shriver SK, Nair HS, Hofstetter R (2013) Social ties and usergenerated content: Evidence from an online social network. *Management Sci.* 59(6):1425–1443.
- Soetevent AR, Kooreman P (2007) A discrete-choice model with social interactions: With an application to high school teen behavior. *J. Appl. Econometrics* 22(3):599–623.
- Stafford TM (2015) What do fishermen tell us that taxi drivers do not? An empirical investigation of labor supply. J. Labor Econom. 33(3):683–710.
- Staiger D, Stock JH (1997) Instrumental variables regression with weak instruments. *Econometrica* 65(3):557–586.
- Su CL, Judd KL (2012) Constrained optimization approaches to estimation of structural models. *Econometrica* 80(5):2213–2230.
- Sun Y, Ishihara M (2018) A computationally efficient fixed point approach to structural estimation of aggregate demand. J. Econometrics 208(2):563–584.
- Toker-Yildiz K, Trivedi M, Choi J, Chang SR (2017) Social interactions and monetary incentives in driving consumer repeat behavior. *J. Marketing Res.* 54(3):364–380.
- van der Lans R, van Bruggen G, Eliashberg J, Wierenga B (2010) A viral branching model for predicting the spread of electronic word of mouth. *Marketing Sci.* 29(2):348–365.
- Vitorino MA (2012) Empirical entry games with complementarities: An application to the shopping center industry. *J. Marketing Res.* 49(2):175–191.
- Wooldridge JM (2010) *Econometric Analysis of Cross Section and Panel Data* (MIT Press, Cambridge, MA).
- Yang J, McAuley J, Leskovec J (2013) Community detection in networks with node attributes. Xiong H, Karypis G, Thuraisingham B, Cook D, Wu X, eds. 2013 IEEE 13th Internat. Conf. Data Mining (IEEE, Piscataway, NJ), 1151–1156.
- Zhu T, Singh V (2009) Spatial competition with endogenous location choices: An application to discount retailing. *Quant. Marketing Econom.* 7(1):1–35.