Estimating Peer Effects in Networked Panel Data

By

Daniel Ian Rock

S.B. Economics
University of Pennsylvania, 2010

SUBMITTED TO THE DEPARTMENT OF MANAGEMENT IN PARTIAL
FULFILLMENT OF THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF SCIENCE IN MANAGEMENT RESEARCH
AT THE
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

JUNE 2016

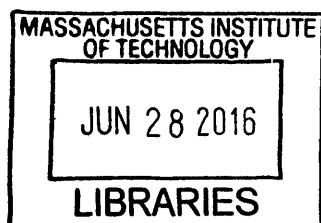Signature of Author:  **Signature redacted**

Department of Management
May 6, 2016

Certified by:  **Signature redacted**

Erik Brynjolfsson
Schussel Family Professor of Management Science
Professor of Information Technology
Thesis Supervisor

Accepted by:  **Signature redacted**

Catherine Tucker
Sloan Distinguished Professor of Management
Professor of Marketing
Chair, MIT Sloan PhD Program

1

# Estimating Peer Effects in Networked Panel Data
## By

Daniel Ian Rock

Submitted to the Department of Management on May 6, 2016 in Partial Fulfillment of the Requirements for the Degree of Master of Science in Management Research

ABSTRACT

After product adoption, consumers make decisions about continued use. These choices can be influenced by peer decisions in networks, but identifying causal peer influence effects is challenging. Correlations in peer behavior may be driven by correlated effects, exogenous consumer and peer characteristics, or endogenous peer effects of behavior (Manski 1993). Extending the work of Bramoullé et al. (2009), we apply proofs of peer effect identification in networks under a set of exogeneity assumptions to the panel data case. With engagement data for Yahoo Go, a mobile application, we use the network topology of application users in an instrumental variables setup to estimate usage peer effects, comparing the performance of a variety of regression models. We find analyses of this type may be especially useful for ruling out endogenous peer effects as a driver of behavior. Omitted variables (especially ones related to network homophily) and violation of the exogeneity assumptions can bias regression coefficients toward finding statistically significant peer effects.

Thesis Supervisor: Erik Brynjolfsson
Title: Schussel Family Professor of Management Science, Professor of Information Technology

## Acknowledgements

**Introduction**

The longstanding econometric challenge of accurately measuring peer effects has recently been the subject of renewed interest (Angrist and Lang 2004; Card and Giuliano 2012; Christakis and Fowler 2007; Dupas et al. 2008; Sacerdote 2001). Recent methodological developments, together with the increased availability of granular data on behavior in networks, offer new opportunities to empirically address questions of social cause and effect. However, the increased number of peer effects studies coincides with considerable discussion of their associated methodological challenges. While the idea that individuals may act under the influence of their peers is widely accepted, measuring and interpreting these effects remains difficult, especially in the context of social networks and digital products.

Research in peer effects often refers to the effects of social influence with respect to isolated diffusion or adoption events. Some of the new datasets available to researchers provide detail on use patterns over time. Online systems in particular offer a wealth of opportunities for observational analysis because the data collection costs are so low. Our aim in this work is to first provide and apply an appropriate model for estimation of peer effects in a dynamic setting, then suggest conditions under which the model is most useful. We discuss some of the econometric issues in using network topology to estimate peer effects. Weak instruments, omitted variables, and violation of the exogeneity assumptions can bias regression coefficients toward finding statistically significant peer effects. Our dataset comes from a Yahoo! mobile browsing application with more than 25 million users. We apply a network-based two-stage least squares estimation procedure to the mobile application's pageview time series at the individual level. Our results suggest that measuring causal peer effects remains tricky. There are many ways to unwittingly overestimate peer effects. In some cases a failure to reject the null of no peer effects is therefore more informative.

Manski (1993) characterizes apparent correlations between individual and peer behaviors to come from three coincident sources: correlated effects, "endogenous" peer effects, and "exogenous" peer effects. In the context of a rainy day, the correlated effect might refer to the tendency of those outside to open umbrellas, endogenous peer effects would be the increased tendency for us to open umbrellas given an observation that others are doing the same, and exogenous peer effects would be the change in our propensity to open an umbrella given an observation of what others are wearing. An ordinary least squares regression of individual umbrella use on peer umbrella use with individual controls would pick up all three effects. Angrist (2014) adds that any regression of an individual outcome on group outcomes alone produces a coefficient of 1. Even if we wanted to study umbrella use on a sunny day, we would still face the issue that exogenous and endogenous effects are not individually identified. Angrist reminds us that here we have Manski's "reflection problem", meaning that "observed behavior is always consistent with the hypothesis that individual behavior reflects mean reference-group behavior". This problem is inescapable without a strategy to distinguish between the different types of peer influence. Moffitt (2001) also discusses issues with identification of peer effects and how they might be solved.

Following the model described by Lee (2007), Bramoullé et al. (2009) offer one way to handle the reflection problem in observational data. Relying on spatial econometrics as a guide, they move beyond groups as traditionally defined toward social network connections to develop reference groups on the individual-level. They provide sufficient conditions for identification of endogenous and exogenous peer effects coefficients in static linear-in-means models. In their model the pre-assigned covariates of nth-degree connections (for n>2 in the basic model) in a social network can be used as instruments for peer behavior. Models of this sort have scarcely been tested in the past. The first purpose of our work is to extend this model to the time-varying case and show that the original identification results still hold for panel data. We then elaborate by investigating challenges to the model's assumptions, offering suggestions for potential applications of this method in business and policy contexts.

Business applications for peer effects estimation have already established a foothold via studies of adoption and diffusion (Aral et al. 2009; Manchanda et al. 2008; Miller and Tucker 2009; Tucker 2008). Scholars in many disciplines have been thinking about these ideas for quite a while (Hartmann et al. 2008). Information about the causes of diffusion is valuable. In the next section we discuss research into adoption and engagement in greater depth. While both "diffusive" behaviors, compared to adoption, engagement is relatively understudied given the business value of proper engagement strategy. Extending the model of Bramoullé et al. to a dynamic panel data application in the style of Arellano and Bond (1991) opens the door for more studies of time-dependent behaviors like engagement and product use. This necessitates additional assumptions to maintain identification of the model parameters. As with similar regression structures, the conditions required for the model to accurately estimate peer effects can pose a challenge.

Identification is a second-order concern if a model can be expected to generate spurious estimates of peer effects. Peer effects two-stage least squares (2SLS) setups (like ours in this paper) and ordinary least squares (OLS) regression estimates might diverge for many reasons, some of them quotidian or mechanical (Angrist 2014). Yet in many cases the relevant standard for an approach to be useful is if the exercise provides actionable information to managers or policy-makers. Knowing, for instance, that a product's engagement is *not* driven by social processes is vital for a marketer considering a viral distribution channel. We find that 2SLS-based network models might help a manager more reliably conclude that their product lacks a social component given the propensity of unobservables to be positively correlated for peers. Much has been made of product virality. But at least in the network we examined the evidence suggests engagement behavior isn't always "contagious".

Section 1 presents key parts of the relevant literature on peer effects and engagement. Section 2 reviews the Bramoullé et al. (2009) model and extends it to the panel data case. Section 3 describes the data and how our dataset was constructed. Section 4 presents estimation results and different model specifications. Section 5 discusses the results and possible extensions. Section 6 concludes.

## Section 1 – Diffusion, Peer Influence, and Estimation of Peer Effects

Social effects can act as an amplifier for policy or business decisions. Human activity is connected; policy-makers or managers should be mindful of the social externalities involved with their choices. Previous work on diffusion and peer influence has spanned many disciplines, each applying a different bundle of methods (Hartmann et al. 2008). Some of that research has focused on generating contagion, or solving the "influence maximization problem". Domingos and Richardson (2001) propose algorithms to maximize lift in marketing efforts via social network interactions. Bakshy et al. (2011) suggest that cascade prediction based on prior influence events in the Twitter network is improved by targeting many influential network nodes. Aral et al. (2013) propose strategies to "engineer" contagion in the presence of homophily, i.e. the tendency of "birds of a feather to flock together" in networks (McPherson et al. 2001). These empirical studies of cascades complement theoretical work discussing or modeling how diffusion might occur (Bass 1969; Granovetter 1978; Jackson and Yariv 2007; Schelling 1971; Watts 2002).

Many studies have also directly measured peer effects at the individual level, under both observational and experimental conditions. For observational datasets, a variety of econometric techniques have been used to estimate social multipliers. Tucker (2008) measures peer effects in the diffusion of a video-messaging technology, comparing the relative size of network externalities for managers and other employees at a bank. This work also used an instrumental variables approach, finding that measures of ego-level importance in the network affected the magnitude of observed adoption externalities. Bollinger and Gillingham (2012) model peer effects in adoption of solar panels using a first differences approach, where lags of peer solar panel installation are used to predict adoption. Consistent estimates of peer effect coefficients in this specification, however, can be problematic because of omitted variables bias even if the "pre-assigned" covariate regressors are truly exogenous to the behavior of interest.

5

Noting that the individual node characteristics and behaviors tend to be correlated with observed networks, Aral et al. (2009) use covariate matching to distinguish between social influences and homophily in adoption of a mobile application. They find that failure to control for homophily leads to overestimation of behavioral contagion effects in adoption by as much as 300-700%. Even matching methods, however, might not fully account for potential bias in estimation. To assess the magnitude of estimation bias, Eckles and Bakshy (2014) compare the results of a randomized experiment to observational techniques. In an experiment with 67 million users on Facebook, naïve observational estimators overstated the experimental measure of contagion by as much as 300%. They point out that some of the difficulty in measuring peer effects comes from the "implausible assumption that the available covariates are sufficient to make peer behavior unconfounded". Thus undercontrolled studies are often more likely to find positive peer effects far in excess of the real externalities.

Social contagion over network ties can operate differentially depending on the individual's "influence" and "susceptibility" (Aral and Walker 2012). This is true beyond the effects of homophily. Network positioning, heterogeneity in tie strength, and dyad-level variation in the sign of social effects combine in real networks to generate the data we observe. Most people can personally relate to the concepts of trendsetters and followers or differences in the strength of social relations. Certainly if peer influence is a fact of life, we have all been either encouraged or discouraged by the actions of others at some point. So-called "chilling effects", when peer behavior slows down the growth rate of others' behaviors, were found in a number of contexts by Goldenberg et al. (2010). They point to adopters of CD players, DVD players, and cellular services waiting for early adopters to act first, creating "hockey stick" shaped growth. Of course growth need not rebound if peer use leads to congestion. The wide range of possibilities for different types of influence underscores that average effects may not tell a sufficiently granular story.

Experimental designs get around the typical worries about exogeneity. Experiments are now common in networks and peer effects research, especially as digital tools have made them cheaper to conduct (Bapna and Umyarov 2012). Some studies have focused on experimentally changing the behavior or product of interest instead of networks or group means. Aral and Walker (2011) create treatment and control groups by randomly embedding viral features of different kinds into a Facebook application and tracking the diffusion. With a hazard model specification, they show a 246% increase in social contagion with the addition of passive message broadcasting. Muchnik et al. (2013) demonstrate herding effects in ratings by manipulating the initial rating of posts on a social news aggregation website. These particular designs escape some of the "perils" of peer effects analysis by separating the subjects of analysis from the treatment (whereas a regression of individual behavior on group mean behavior does not). One concern for networked experiments is interference, or the tendency for the stable unit treatment value assumption (SUTVA) to fail when treatments can diffuse into control groups. Athey et al. (2015) discusses p-value calculation in networked contexts.

Other experimental designs focus on changing network exposure conditions and observing differences in diffused behaviors. This might entail changing the connections in a network or experimentally altering the information moving between connected individuals. In the Facebook network, for example, people who were exposed to signals about friends' information (in the form of different urls) were found more likely to transmit that information themselves (Bakshy et al. 2012). In that study it was the exposure condition that the researchers randomly assigned. In contrast, the treatment in Bapna and Umyarov (2012) was being exposed to a peer who had adopted a paid service. This later design is more analogous to the observational models commonly applied to the study of peer effects.

Many of these studies focus on one-shot behaviors like adoption, ratings, or clicks. But actions that vary over time may often carry more value. Fader et al. (2005) point out that measuring customer lifetime value (CLV), for example, may come from individual-level conditional expectations or aggregate-level sales projections. They offer a taxonomy of probability models for estimating CLV, noting that "recency, frequency, and monetary value" or "RFM" of transactions can be used to estimate "residual lifetime value" (customer value unexplained by covariates). Though one of many possible frameworks, but RFM is conceptually closer to engagement and use than it is to adoption. Peer effects in usage may differ from peer effects in adoption, with large implications for businesses looking to monetize user behavior. Companies

can set usage-based pricing structures and subsidize influential individuals ("opinion leaders") to generate value or stimulate diffusion (Ghose and Han 2011; Iyengar et al. 2011). These studies on engagement rely on observational data as well. As with Bollinger and Gillingham's work on solar panel adoption, they use fixed effects in an attempt to handle the reflection problem. Aral and Nicolaides (2016) examine peer effects in exercise over time, relying on exogeneity of the weather conditions for peers in other cities in a continuous IV setup. This research design more convincingly argues for "contagion" in healthy behaviors than observational data would. Future work might compare the parameter estimates from studies of this sort to those retrieved from observational network methods.

Whether an observational or experimental study design, it isn't clear that "the coefficient on group averages in a multivariate model of endogenous peer effects reveals the action of social forces" (Angrist 2014; Shalizi and Thomas 2011). Consider, as Angrist does, a setup where we estimate the OLS and instrumental variables (IV) coefficients in a regression of individual on group behavior (group average behavior is the instrumented variable). The peer effect is represented as the divergence between the OLS and IV of the relationship between covariates and outcomes, but these estimands can differ for reasons other than peer effects. Many standard econometric problem apply here, plus a few more unique to networks. In a many weak instruments situation, the IV estimates will be biased toward OLS estimates. "Common variance components in outcomes" can produce correlations that resemble social effects, but are driven by omitted factors. And deciding which observation of "the network" is the true latent representation of links between units can be a judgment call.

The model we develop is appropriate for observational networked data where outcome behavior and at least one of the covariates vary over time. It is principally based upon the work of Bramoullé et al. (2009), which explicitly leverages network topology to create local group means of outcome and explanatory variables. The presence and linearly independent arrangement of intransitive triads of individuals in the network permits identification results to hold. In other words, a network where all friends-of-friends are also friends will lead to degenerate estimation outcomes. That kind of model reduces to a regression of outcomes on leave-out group means and covariates.

## Section 2 – Extension of Bramoullé, Djebbari, and Fortin's Social Network Model (BDF)

*A review of the static network model*

The baseline structural model used describes the (linear) relationship between the individual behavior vector $\mathbf{y}$, the vector of connected peers' behaviors $\mathbf{Gy}$, the vectors of exogenous (i.e. pre-assigned in this case) covariates $\mathbf{x}$, and the vectors of pre-assigned covariates for connected peers $\mathbf{Gx}$. $\mathbf{G}$ is the row-normalized adjacency matrix representing the social network of peers. Entries in $\mathbf{G}$ are zero if two individuals are not connected ($G_{ij} = 0$ if there is no i,j dyad). For simple averaging of peer behavior, entries in each row (representing the peers of individual $i$) will sum to 1 and will have the value $1/n$ where n is the count of nonzero entries in the row. This is equivalent to a simple leave-out mean. To review, following the simplest version of Bramoullé et al. (2009) (or "BDF" for the rest of the paper) we have the equation:

$$y = \alpha\iota + \beta Gy + \gamma x + \delta Gx + \epsilon; \quad \mathbb{E}[\epsilon \mid x] = 0 \quad (1)$$

With $|\beta| < 1$ and $(I - \beta G)^{-1} = \sum_{k=0}^{\infty} \beta^k G^k$, we have that:

$$y = \alpha(I - \beta G)^{-1}\iota + (I - \beta G)^{-1}(\gamma I + \delta G)x + (I - \beta G)^{-1}\epsilon \quad (2)$$

$$y = \left(\frac{\alpha}{1 - \beta}\right)\iota + \gamma x + (\gamma\beta + \delta)\sum_{k=0}^{\infty} \beta^k G^{k+1}x + \sum_{k=0}^{\infty} \beta^k G^k \epsilon \quad (3)$$

Equations (2) and (3) describe the reduced form structure of the model where $y$ is purely a function of observables $\mathbf{x}$ and network $\mathbf{G}$. We then can write the expected average peer behavior as:

$$\mathbb{E}(Gy|x) = \frac{\alpha}{(1-\beta)}\iota + \gamma Gx + (\gamma\beta + \delta)\sum_{k=0}^{\infty}\beta^k G^{k+2}x \quad (4)$$

BDF goes on to show in their proposition 1 that if $\gamma\beta + \delta \neq 0$ and the matrices $\mathbf{I}$, $\mathbf{G}$, and $\mathbf{G^2}$ are linearly independent, then social effects are identified. Additionally if those 3 matrices are linearly dependent and no individual is isolated, social effects are not identified (the paper has more detailed results). Without linear independence, then it is not possible to find an identifying instrument for $\mathbf{Gy}$ in equation (1) above. Otherwise we can use network structure to find exogenous covariates that instrument for $\mathbf{Gy}$. This is to say that friends-of-friends' covariates ($\mathbf{G^2x}$) or even deeper network connections' covariates ($\mathbf{G^n x}$ for n>2) can be used to predict friends' behavior in an IV setup. While for our purposes we will restrict the model in this analysis to one connected component of a network, writing a block diagonal matrix $\boldsymbol{\Gamma}$ where the diagonal is composed of disjoint subnetworks $\mathbf{G_1},...,\mathbf{G_n}$ is a natural extension of this model. Including network fixed effects, BDF generalizes the model to partially deal with correlated effects. Note that this does not deal with the tendency for networks to be formed *because* of behaviors. We still must assume that the network and network formation process is strictly exogenous to vector y. With the correlated effects setup, for a network $k$:

$$y_k = \alpha_k + \beta G_k y_k + \gamma x_k + \delta G_k x_k + \epsilon_k; \quad \mathbb{E}[\epsilon_k \mid \alpha_k, G_k, x_k] = 0 \quad (5)$$

Subtracting the local network average from the individual, BDF shows the result that $\mathbf{I}$, $\mathbf{G}$, $\mathbf{G^2}$, and $\mathbf{G^3}$ must be linearly independent in addition to the assumption that there are social effects, i.e. $\gamma\beta + \delta \neq 0$. In this case the valid identifying instrument set will be $(I - G_k)G_k^2 x_k, (I - G_k)G_k^3 x_k, ..., (I - G_k)G_k^n x_k$. This "within local transformation" is analogous to a first difference estimator where the fixed effect for network k is differenced out across the network. It is this model and Arellano and Bond (1991) which inspire the panel model. We present the setup for a single connected network, though the results can be extended to include multiple networks with correlated effects in a stacked matrix.

*The Simple Panel Model*
For a single connected network, if t indexes time, we have:

$$y_t = \alpha\iota + \beta G_t y_t + \gamma x_t + \delta G_t x_t + \epsilon_t; \quad \mathbb{E}[\epsilon_t \mid x_t, G_t] = 0 \ \forall t \quad (5)$$

And analogously to (2) and (3), with $|\beta| < 1$, $I - \beta G_t$ and $(I - \beta G_t)^{-1} = \sum_{k=0}^{\infty}\beta^k G_t^k$, we have that:

$$y_t = \alpha(I - \beta G_t)^{-1}\iota + (I - \beta G_t)^{-1}(\gamma I + \delta G_t)x_t + (I - \beta G_t)^{-1}\epsilon_t \quad (6)$$

$$y_t = \left(\frac{\alpha}{1-\beta}\right)\iota + \gamma x_t + (\gamma\beta + \delta)\sum_{k=0}^{\infty}\beta^k G_t^{k+1}x + \sum_{k=0}^{\infty}\beta^k G_t^k\epsilon_t \quad (7)$$

$$\mathbb{E}(G_t y_t | x_t) = \frac{\alpha}{(1-\beta)}\iota + \gamma G_t x_t + (\gamma\beta + \delta)\sum_{k=0}^{\infty}\beta^k G_t^{k+2}x_t \quad (8)$$

(see Appendix A.1 in Bramoullé et al. (2009) for a proof that proves results for all models once they have been reduced to equations analogous to (5)-(8)).

*The Panel Model with Fixed Effects and Static Networks*
This setup has no means of handling the propensity for "common shocks" to affect networks over time or within connected components. If we want to model time-specific fixed effects, we have a structural equation of the form:

$$y_t = \lambda_t \iota + \alpha\iota + \beta G_t y_t + \gamma x_t + \delta G_t x_t + \epsilon_t \qquad (9)$$

Using first differences, this setup too can be used to derive the reduced form in analogous way. Let

$$\Delta y_t = y_{t+1} - y_t, \ \Delta x_t = x_{t+1} - x_t, \ \Delta\epsilon_t = \epsilon_{t+1} - \epsilon_t, \ \Delta\lambda_t = \lambda_{t+1} - \lambda_t$$

By differencing two equations at time t and time t+1, we have:

$$\Delta y_t = \Delta\lambda_t\iota + \beta(G_{t+1}y_{t+1} - G_t y_t) + \gamma\Delta x_t + \delta(G_{t+1}x_{t+1} - G_t x_t) + \Delta\epsilon_t;$$
$$\mathbb{E}[\Delta\epsilon_t \mid \Delta x_t, G_{t+1}, G_t] = 0 \ \forall t \quad (10)$$

Assume for now that the latent network is stable and $G_t$ is isomorphic for all time t. We'll denote this graph $G_f$ for "final graph". In this case the problem reduces cleanly as we can factor the final graph out. We get:

$$\Delta y_t = \Delta\lambda_t\iota + \beta G_f \Delta y_t + \gamma\Delta x_t + \delta G_f \Delta x_t + \Delta\epsilon_t \quad (11)$$

With reduced form, by an argument similar to the ones above:

$$\Delta y_t = \frac{\Delta\lambda_t}{1-\beta}\iota + \gamma\Delta x_t + (\gamma\beta + \delta)\sum_{k=0}^{\infty}\beta^k G_f^{k+1}\Delta x_t + \sum_{k=0}^{\infty}\beta^k G_f^k \Delta\epsilon_t \quad (12)$$

$$\mathbb{E}(G_f \Delta y_t \mid \Delta x_t) = \frac{\Delta\lambda_t}{1-\beta}\iota + \gamma G_f x_t + (\gamma\beta + \delta)\sum_{k=0}^{\infty}\beta^k G_f^{k+2}\Delta x_t \quad (13)$$

Incidentally this will be the reduced form equation that applies to the Yahoo! Go network. We can only observe the network on the final day of a 28 day period, and therefore take the network at the end of the period as the latent "true" network (this assumes that if two people were friends by the end of the month, they were likely friends at the beginning of the month too). Where these results fail, however, is in the case that $G_t \neq G_f$ for all time t. In that case equation 10 does not reduce as cleanly. It may be that the network is measured with error; some observed connections might not be active while other unobserved connections are activated. We assume we can observe the true latent network, though previous work has developed methods to more concretely describe links between nodes (A. Goldenberg et al. 2010).

*The Panel Model with Fixed Effects and Dynamic Networks*

Returning to (10), to use the series expansion shortcut to the reduced form we must find some matrix which represents the change in the network over time. Let us now modify **G** slightly to use matrix **F**, which will represent a version of **G** that is not yet row-normalized. $F_{ij} = 1$ if individual i and individual j are connected. **Fy** is then the sum of peer behaviors for a given individual. We have:

$$\Delta y_t = \Delta\lambda_t\iota + \beta(F_{t+1}y_{t+1} - F_t y_t) + \gamma\Delta x_t + \delta(F_{t+1}x_{t+1} - F_t x_t) + \Delta\epsilon_t;$$
$$\mathbb{E}[\Delta\epsilon_t \mid \Delta x_t, \Delta F_t] = 0 \ \forall t \quad (14)$$

so it follows that:

$$\Delta y_t - \beta(F_{t+1}y_{t+1} - F_t y_t) = \Delta\lambda_t\iota + \gamma\Delta x_t + \delta(F_{t+1}x_{t+1} - F_t x_t) + \Delta\epsilon_t$$

The researcher must decide on the precise timing of peer effects to allow for separation of changing network effects from changing behavior effects. We have two boundary conditions. In the "early" condition, former friends are discarded and the relevant adjacency matrix for the next period will be $F_{t+1}$. In the "late" condition, new friends have no peer influence and the relevant adjacency matrix is $F_t$. In general, we can define a parameter (or potentially a vector of parameters) $\rho \in [0,1]$ to determine relative influence levels

9

of new friends and old friends. This weighting value determines which convex combination of friends to use.

$$F_{\rho,t} = \rho F_{t+1} + (1 - \rho)F_t \quad (15)$$

With $\varrho$ selected, we can renormalize $F_{\varrho,t}$ row-wise to get $G_{\varrho,t}$, and then we have:

$$\Delta y_t - \beta G_{\rho,t}\Delta y_t = \Delta\lambda_t \iota + \gamma\Delta x_t + \delta G_{\rho,t}\Delta x_t + \Delta\epsilon_t \quad (16)$$

And then using the representation of $|\beta| < 1$, $(I - \beta G_\rho)^{-1} = \sum_{k=0}^{\infty}\beta^k G_\rho^k$, equation (12) will hold subbing in $G_{\varrho,t}$ for $G_t$. Appendix A briefly discusses a few ways to calculate $G_\varrho$. For the "between" estimator, where we have individual specific fixed effects, there is no problem factoring out the network after de-meaning the vectors in the equation. The network is stable in cross-section after fixing the time period.

*General features of all of the models*

All of these models struggle to handle some of the key problems with measuring peer effects. Nowhere do we fix the problem of distinguishing genuine social contagion from homophily. Correlated effects can be differenced out at the network-level, but subgraphs may contain their own localized correlated effects that further influence connected peers. To the extent that pre-assigned characteristics of peers influence behaviors (without, somehow, influencing formation of ties), we can calculate an estimate of an average peer effect of behavior contingent on observable types of homophily. Yet rarely do researchers have sufficient data on covariates to control for every type of correlation in tie formation. Additionally, there is often good reason to think that a given behavior of interest generates peer connections from period to period. In that case the moment restrictions do not hold. It is therefore difficult to fully trust a causal interpretation of the peer effects coefficients generated by these models. The identification of these parameters is subordinate to the concern that they are fragile. In other words, we should maintain skepticism of the estimated parameters because previous research on behavior in networks suggests these models will often fail to fully describe social dynamics.

Pessimism aside, the structure of the reduced form of the type in equations (3), (7), and (12) permits a robustness check on the assumptions of the model. We could run a kind of truncated reduced form regression of outcomes on covariates, friends' covariates, friends-of-friends' covariates,…, n-distant friends' covariates, etc. Since $|\beta| < 1$, we would expect to see coefficients on friends' covariates declining exponentially in network distance. This coefficients would of course be biased (likely upward) by the omitted variables calculated by letting n go to infinity. Nevertheless, normalizing the coefficient on the immediate neighbors' covariates to be 1, we should see coefficients rapidly shrink as a function of network distance. This kind of check would be most useful in the case that the model is just identified. Otherwise a test of overidentification suffices. So far we have also assumed that social effects exist. In the case that they *do not* exist, we expect that the first stage of the IV regression collapses. This condition of no social effects implies that the 2SLS estimates will not differ greatly from OLS estimates in expectation. That weak instrument scenario presents an asset: either the included covariates have no social content, the behavior of interest is unlikely to be contagious, or both. These are useful facts for managers to know.

**Section 3 – The Yahoo! Go Data and Network**

Our network, engagement, and covariate data come from Yahoo! Go, a mobile application designed for browsing behavior. Yahoo! designed Go as a way for online users to access Yahoo! services on their mobile phones and/or PDAs. Launched in July 2007 and discontinued in January 2010, the network of users had over 27 million members. Go services enabled users to check sports scores, look up stock quotes, send and receive email, search, and read news (see Appendix C for an image of one of the application screens). Each user has a unique Yahoo! ID across all Yahoo! services (mobile or otherwise). These IDs

have been anonymized by the company and all users under age 18 have been removed from the dataset. Go, however, does not support instant messaging and social activities like link sharing, commenting, or liking were not included on the Go! Platform. Go activities should not be expected to exhibit network externalities in usage. This is to say that though adoption might spread via word-of-mouth or similar channels, use is *not* likely to be social on the basis that an observable network exists.

The exogeneity of the network with respect to the behavior of interest is assumed in the model setup we have described. In practice these types of networks are hard to come by, especially with the proliferation of social or "viral" application features in digital products. The separation of network formation from the behavior of interest makes the Go network an attractive context in which to apply the model. We wish to rule out the criticism that our results might be driven by network endogeneity in discussing the performance and estimation results of our model specifications. Go serves that purpose, facilitating a more nuanced discussion of the research value of network topology-based panel models.

The network of interest is formed by instant messaging (IM) behavior, and we have daily use data in the form of mobile page views for October 2007. We also have covariate data on user demographics, including age, gender, and friend counts (degree). Using the IM connections, we also build a series of network covariates which are static over the course of our panel (since we use the end-of-month network to represent the latent network). We use the total mobile pageviews because the desktop behaviors might also be linked to network formation processes. Starting with the 27 million Yahoo! Messenger users, we narrow the user set to adopters of Go who also have friends using Go (21,896 users), and then perform the analysis on the largest connected component to satisfy the condition that $\mathbf{I}$, $\mathbf{G}$, and $\mathbf{G^2}$ are linearly independent. This leaves us with 2,203 users with use behavior observed over 28 days. Figure 1.A depicts the network of users, while figure 1.B zooms in on the largest connected component. Figure 1.C shows the degree distribution of nodes in the largest connected component and 1.D shows the distribution of node counts in all connected components of users *except* the largest one. Our largest connected component is exceptionally large in comparison to the others.



Figure 1.A – User Network (21,896 users)



Figure 1.B – Largest Connected Component (2,203 users)

11

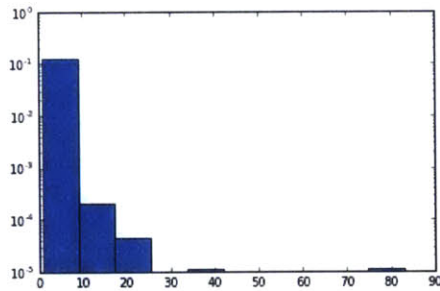Figure 1.C – Degree Distribution (logged frequencies) in largest component



Figure 1.D – Size Distribution of Other Connected Components (11,108 in total)

In our largest connected component, the degree distribution closely follows an exponential density and the network diameter is 21. This "scale-free" characteristic is common in networks (Barabási and Albert 1999). Aside from the largest connected component, there are 11,108 smaller networks of users which either fail to have a large enough user pool to develop precise estimates, fail to meet the identification conditions of the model, or both. Summary statistics for the largest connected component are reported below (Table 1).

| Table 1 – Summary Statistics | Mean | Standard Deviation | Count |
|---|---|---|---|
| Age | 35.53 | 7.93 | 2203 |
| Friend Age | 35.25 | 4.51 | 2203 |
| Friend-of-Friend Age | 35.36 | 2.94 | 2203 |
| Gender | 0.77 | 0.43 | 2203 |
| Friend Gender | 0.76 | 0.28 | 2203 |
| Friend-of-Friend Gender | 0.76 | 0.18 | 2203 |
| Degree | 6.41 | 5.84 | 2203 |
| Friend Degree | 10.07 | 5.17 | 2203 |
| Friend-of-Friend Degree | 9.83 | 3.68 | 2203 |
| PageRank | 0.000046 | 0.000029 | 2203 |
| Friend PageRank | 0.000067 | 0.000022 | 2203 |
| Friend-of-Friend PageRank | 0.000064 | 0.000013 | 2203 |
| Use (pageviews) | 16.90 | 44.02 | 61684 |
| Lagged Use (1 Day) | 17.09 | 44.40 | 61684 |
| Lagged Use (2 Days) | 17.16 | 44.55 | 61684 |
| Friend Use | 18.70 | 27.39 | 61684 |
| Lagged Friend Use (1 Day) | 18.91 | 27.62 | 61684 |
| Lagged Friend Use (2 Day) | 19.03 | 27.77 | 61684 |
| Friend-of-Friend Use | 18.45 | 14.72 | 61684 |
| Lagged Friend-of-Friend Use (1 Day) | 18.66 | 14.79 | 61684 |
| Lagged Friend-of-Friend Use (2 Day) | 18.76 | 14.83 | 61684 |

Table 1 – Summary Statistics for Selected Variables

It is not surprising that friend and friend-of-friend summary statistics are so similar to individual statistics (everyone is someone's friend). We do see a common feature of human networks that the average degree of friends is higher than the individual's degree. This is because the most connected individuals are more likely to be friends with many people. We also calculate the PageRank score for all nodes in the network as a means of measuring network centrality (Page et al. 1998). This will be used as a control in our peer effects regressions. The average age is 35.53 with a standard deviation of 7.93. Gender is coded as 1 for men and 0 for women. Since we only have use (mobile pageviews) as a time-varying characteristic of individuals in our network, we will use lags of use as our pre-assigned covariate for individuals, friends, and friends-of-friends. In the next section we present the regression results from estimation of our panel specifications.

## Section 4 – Regression Results

The tables are organized as follows: Tables 2A and 2B show estimates from OLS regressions of individual use (in mobile pageviews) on the average friend use and other covariates. Tables 3A and 3B show the IV estimates for the same equations, using all available friend-of-friend covariates as instruments for the endogenous average friend use. For example, we have in columns 3 and 4 of Table 2A pooled and fixed effect versions (respectively) of a regression of individual use on friends' use, lagged individual use (1 and 2 days lagged) and lagged friends' use. We use lagged friends-of-friends' use (1 and 2 day) to

instrument for friends' use in columns 1 and 2 of Table 3A. When we include more demographic covariates to predict individual use, we similarly use friends-of-friends covariates of the same type as instruments. First stage results are presented in Tables 4A and 4B in Appendix B. Tables 5A-C present Sargan J statistics and Cragg-Donald Wald F statistics and Stock and Yogo (2005) weak identification critical values for the 6 IV specifications.

The "full model" containing all covariates and lagged usage behaviors for the individual and their friends has the following equation (similar to equation 9 in section 2):

$$Use_t = \alpha + \beta G_t Use_t + \gamma_1 LagUse1_t + \gamma_2 LagUse2_t + \gamma_3 Age + \gamma_4 Gender + \gamma_5 Degree$$
$$+ \gamma_6 PageRank + \delta_1 G_t LagUse1_t + \delta_2 G_t LagUse2_t + \delta_3 G_t Age + \delta_4 G_t Gender$$
$$+ \delta_5 G_t Degree + \delta_6 G_t PageRank + \lambda_t + \epsilon_{i,i}$$

We substitute in the network we observe at the end of the period of analysis $G_f$ (in row-normalized form) for $G_t$ for each day in the 28 day period we examined. The other models use a subset of the variables included in this larger equation. Multiplying a vector by $G$ therefore returns a vector containing the average for an individual's friends' values (use or covariate). We use 1 and 2 day lags in these specifications.

| Table 2A – OLS Dep. Var.: Use (Pageviews) | (1) OLS - Pooled | (2) OLS - Fixed Effects | (3) OLS - Pooled with Lags | (4) OLS - Fixed Effects with Lags |
|---|---|---|---|---|
| Friend Use | 0.0665*** | 0.0560*** | 0.0444*** | 0.0299*** |
|  | (0.00636) | (0.00635) | (0.00697) | (0.00696) |
| Lag Use (1 Day) |  |  | 0.299*** | 0.298*** |
|  |  |  | (0.0122) | (0.0121) |
| Lag Use (2 Day) |  |  | 0.189*** | 0.191*** |
|  |  |  | (0.00986) | (0.00990) |
| Friend Lag Use (1 Day) |  |  | 0.00295 | 0.00363 |
|  |  |  | (0.00766) | (0.00772) |
| Friend Lag Use (2 Day) |  |  | -0.000290 | 0.00975 |
|  |  |  | (0.00648) | (0.00653) |
| Constant | 15.65*** | 16.04*** | 7.557*** | 2.844*** |
|  | (0.205) | (0.941) | (0.226) | (0.672) |
| Observations | 61,684 | 61,684 | 57,278 | 57,278 |
| R-squared | 0.002 | 0.005 | 0.180 | 0.184 |
| Time FE | NO | YES | NO | YES |

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

**Table 2A – OLS Estimates**

14

| Table 2B - OLS Dep. Var.: Use (Pageviews) | (1) OLS - Pooled with Demographics | (2) OLS - FE with Demographics | (3) OLS - Pooled Full Model | (4) OLS - FE Full Model |
|---|---|---|---|---|
| Friend Use | 0.0685*** | 0.0580*** | 0.0431*** | 0.0286*** |
| | (0.00640) | (0.00638) | (0.00695) | (0.00695) |
| Lag Use (1 Day) | | | 0.297*** | 0.296*** |
| | | | (0.0121) | (0.0121) |
| Lag Use (2 Day) | | | 0.187*** | 0.189*** |
| | | | (0.00983) | (0.00988) |
| Friend Lag Use (1 Day) | | | 0.00208 | 0.00268 |
| | | | (0.00771) | (0.00777) |
| Friend Lag Use (2 Day) | | | -0.00134 | 0.00863 |
| | | | (0.00650) | (0.00654) |
| Age | 0.00735 | 0.00828 | -0.00574 | -0.00554 |
| | (0.0218) | (0.0218) | (0.0192) | (0.0192) |
| Gender | 3.125*** | 3.087*** | 1.784*** | 1.769*** |
| | (0.453) | (0.453) | (0.428) | (0.428) |
| Friend Age | 0.116*** | 0.119*** | 0.0522 | 0.0525 |
| | (0.0411) | (0.0411) | (0.0382) | (0.0383) |
| Friend Gender | -3.558*** | -3.535*** | -1.243** | -1.230** |
| | (0.580) | (0.580) | (0.560) | (0.559) |
| Degree | | | 0.344*** | 0.343*** |
| | | | (0.106) | (0.106) |
| PageRank | | | -21,569 | -21,461 |
| | | | (19,150) | (19,124) |
| Friend Degree | | | 0.0745 | 0.0774 |
| | | | (0.0851) | (0.0850) |
| Friend PageRank | | | -12,379 | -12,515 |
| | | | (16,789) | (16,757) |
| Constant | 11.57*** | 11.84*** | 4.496*** | -0.255 |
| | (1.513) | (1.773) | (1.555) | (1.685) |
| Observations | 61,684 | 61,684 | 57,278 | 57,278 |
| R-squared | 0.003 | 0.006 | 0.182 | 0.185 |
| Time FE | NO | YES | NO | YES |

Robust standard errors in parentheses

*** $p<0.01$, ** $p<0.05$, * $p<0.1$

**Table 2B – OLS Estimates**

| Table 3A – IV Dep. Var.: Use (Pageviews) Instrumented: Friend Use | (1) IV - Pooled with Lags | (2) IV - FE with Lags | (3) IV - Pooled with Demographics | (4) IV - FE with Demographics |
|---|---|---|---|---|
| Friend Use | 0.618 | 0.536* | 0.722*** | 0.633*** |
| | (0.440) | (0.315) | (0.253) | (0.220) |
| Lag Use (1 Day) | 0.295*** | 0.295*** | 0.294*** | 0.294*** |
| | (0.0131) | (0.0125) | (0.0126) | (0.0124) |
| Lag Use (2 Day) | 0.187*** | 0.188*** | 0.186*** | 0.187*** |
| | (0.00986) | (0.0101) | (0.00983) | (0.0100) |
| Friend Lag Use (1 Day) | -0.178 | -0.154 | -0.209** | -0.182*** |
| | (0.141) | (0.0995) | (0.0818) | (0.0698) |
| Friend Lag Use (2 Day) | -0.110 | -0.0914 | -0.127*** | -0.109** |
| | (0.0839) | (0.0632) | (0.0483) | (0.0442) |
| Age | | | -0.0377 | -0.0342 |
| | | | (0.0247) | (0.0235) |
| Gender | | | 2.637*** | 2.489*** |
| | | | (0.671) | (0.622) |
| Friend Age | | | -0.0120 | -0.00357 |
| | | | (0.0507) | (0.0481) |
| Friend Gender | | | -2.487*** | -2.405*** |
| | | | (0.664) | (0.638) |
| Constant | 2.544 | 3.488 | 3.200 | 3.702 |
| | (3.805) | (3.328) | (2.003) | (2.266) |
| Observations | 57,278 | 57,278 | 57,278 | 57,278 |
| R-squared | 0.079 | 0.106 | 0.040 | 0.075 |
| Time FE | NO | YES | NO | Yes |

Robust standard errors in parentheses
*** $p<0.01$, ** $p<0.05$, * $p<0.1$

**Table 3A – IV Estimates**

16

| Table 3B – IV Dep. Var.: Use (Pageviews) Instrumented: Friend Use | (1) IV - Pooled Full Model | (2) IV - FE Full Model |
|---|---|---|
| Friend Use | 0.758** | 0.555* |
|  | (0.381) | (0.330) |
| Lag Use (1 Day) | 0.293*** | 0.293*** |
|  | (0.0127) | (0.0124) |
| Lag Use (2 Day) | 0.185*** | 0.186*** |
|  | (0.00979) | (0.0101) |
| Friend Lag Use (1 Day) | -0.218* | -0.157 |
|  | (0.119) | (0.101) |
| Friend Lag Use (2 Day) | -0.132* | -0.0924 |
|  | (0.0697) | (0.0637) |
| Age | -0.0305 | -0.0236 |
|  | (0.0265) | (0.0242) |
| Gender | 2.844*** | 2.546*** |
|  | (0.783) | (0.711) |
| Friend Age | -0.0267 | -0.00474 |
|  | (0.0617) | (0.0557) |
| Friend Gender | -2.349*** | -2.039** |
|  | (0.876) | (0.799) |
| Degree | 0.377*** | 0.367*** |
|  | (0.116) | (0.111) |
| PageRank | -24,778 | -23,855 |
|  | (20,901) | (19,964) |
| Friend Degree | -0.228 | -0.143 |
|  | (0.184) | (0.161) |
| Friend PageRank | 31.73 | -3,450 |
|  | (19,717) | (18,370) |
| Constant | 3.849** | 4.262** |
|  | (1.803) | (2.049) |
|  |  |  |
| Observations | 57,278 | 57,278 |
| R-squared | 0.025 | 0.102 |
| Time FE | NO | YES |

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

**Table 3B – IV Estimates**

| Table 5A - Over ID and Weak ID Tests | Pooled with Lags | FE with Lags |
|---|---|---|
| J Statistic | 0.974 | 0.12 |
| J Statistic P-Value | 0.3237 | 0.729 |
| Cragg-Donald Wald F Statistic | 16.969 | 26.749 |
| Stock-Yogo Weak ID 10% Size Critical Value | 19.93 | 19.93 |
| Table 5B - Over ID and Weak ID Tests | Pooled with Demographics | FE with Demographics |
| J Statistic | 21.339 | 21.475 |
| J Statistic P-Value | 0.0001 | 0.0001 |
| Cragg-Donald Wald F Statistic | 18.99 | 23.21 |
| Stock-Yogo Weak ID 10% Size Critical Value | 10.27 | 10.27 |
| Table 5C - Over ID and Weak ID Tests | Pooled Full Model | FE Full Model |
| J Statistic | 26.344 | 27.794 |
| J Statistic P-Value | 0.0001 | 0 |
| Cragg-Donald Wald F Statistic | 34.187 | 6.53 |
| Stock-Yogo Weak ID 10% Size Critical Value | 11.12 | 11.12 |

The OLS regression of individual pageviews on peer average pageviews would suggest there are peer effects in our network. A simple regression of individual on peer use without controls suggests a 1 view increase in the peer average corresponds to a 0.0665 view increase for individuals (Table 2A column 1). With more controls, this increase remains statistically significant but declines to 0.0299 once we include fixed effects for the day of the month and lagged use for the individual and peers. Unsurprisingly, lagged individual use is the best predictor of future individual use. An extra pageview one day prior and two days prior is associated with a respective (statistically significant at 5%) increases of around .3 and .19 pageviews in nearly all models. The full OLS model (Table 2B column 4) with fixed effects predicts a statistically significant 0.0286 extra pageviews for an increase in mean peer pageviews of 1. Gender is significant and also predicts more engagement (1.77 additional pageviews if the individual is male), though having more male friends predicts a decline in individual engagement. Individuals of higher degree also tend to have an additional 0.34 pageviews for every additional friend they have. The purely correlational analysis appears consistent with peer effects, though these correlations are misleading. We should note here that aside from use and gender of friends, there are no statistically significant friend-based measures in the full model. This is the first indication that certain types of peer covariates may not have exogenous peer effects (to use Manski's framing).

The IV regressions offer mixed evidence on the existence of peer effects. In Table 3A column 1, we see the results for the pooled IV with only lagged individual and friend behavior (no fixed effects for the day). Friend Use is instrumented by friend-of-friend lagged use of 1 and 2 days. Once again we have the statistically significant lagged individual coefficients of about 0.3 and 0.19 for 1 and 2 day lags respectively. But the average friend pageviews are not significantly different from 0. Adding time fixed effects (column 2) gets friend use to marginal but not convincing significance (10%). The specifications in columns 3 and 4 add age and gender covariates for the individual and friends, instrumenting for friend use with the lagged friend-of-friend use as before in addition to friend-of-friend "average" gender and age. Now peer effects appear to be statistically significant at the 1% level and quite strong! In column 4, a 1 pageview increase in peer averages appears to cause a 0.633 pageview increase in the individual's use. The full specification (Table 3B column 4) tempers the apparent effect: here an increase in peer average pageviews of 1 corresponds to a 0.555 increase in individual views at a 10% significance level. The coefficients on friends' covariates are quite close to zero. This is an important red flag to which we will return soon.

## Discussion of Results

How can we reconcile the mixed indications for evidence of peer effects? A combination of the overidentification and weak instruments tests is informative. Looking at Table 5A, the J statistics are small

18

enough that we fail to reject at even a 10% significance level the null hypothesis that the instruments are exogenous. Under the assumption that contemporaneous peer pageviews are indeed the only endogenous variable, lags of friends-of-friends' use behavior appear to be valid instruments. Yet once we add in instruments that reflect demographics (age and gender), Table 5B tells us that at least one of the instruments is not exogenous. We reject the null hypothesis for the Sargan Overidentification Test at nearly any reasonable significance level (this is true for the full model in Table 5C as well). The exclusion restriction fails for at least one covariate. The correlation with the error term of friend-of-friend demographics suggests an unobserved covariate related to the individual's outcomes. It would appear that the evidence for peer effects is relatively weak. Either we fail to include covariates which are likely related to the data-generating process, or our peer effects estimates are too imprecise to be convincing.

We assumed when applied this model BDF's "natural" condition that $\gamma\beta + \delta \neq 0$, i.e. there are social effects in this context. Maybe there are social effects, but evidently the friend-of-friend covariate instruments which make that condition true are in fact endogenous. These are the unobservable "shared components of variance in outcomes" which make precisely measuring endogenous peer effects challenging. Models relying on network structure to generate exogenous network instruments fall victim to the likely condition of widespread homophily. Setups such as the one we have presented and the original static version in BDF are especially prone to find positive peer effects whenever positive sorting on unobserved covariates is correlated with the behavior of interest.

If "birds of a feather flock together", then they likely roost together as well (as it were). Indeed models like these effectively stack the deck in favor of finding peer effects by first isolating connected individuals and only then applying an estimation technique within similar groups. Regressions on group leave-out average outcomes are also problematic, but might have additional noise in the definition of "peer" so as to avoid intensified homophily. Furthermore in some applications it may be a problem that the matrix $G^n$ represents paths of length n, including paths that return to the individual (though not if assumptions are strictly satisfied). This can create mechanical nonzero correlations between instrumented peer behavior and outcomes that do not exist for "leave-out" matrices. It is also possible that there are no exogenous or endogenous peer effects, in which case any coefficients on peers we have found is due to statistical noise. We might have weak instruments. Our first stage estimates (Appendix B) have predictive power though. So why isn't there substantial evidence for peer effects? Yahoo! Go lacks social features, but the network exhibits homophily. It would be surprising if we found endogenous peer effects. The Yahoo! Go product engagement is driven primarily by previous engagement. Marketing strategies to leverage "virality" in behaviors for Go would likely fail. A loyal user base (however small) existed, but the lack of a social component likely contributed to the decision to stop support for Go.

## Conclusion

Given the multitude of reasons that a 2SLS IV estimator and an OLS estimator might diverge, of which peer effects is one of many, what is the place for network structure-based IV setups like the one we have used and BDF? Because the model setup is so fragile, estimation results reflecting social multipliers are unlikely to be trustworthy. Researchers must carefully consider whether using network instruments is appropriate given the context. But arguably there is a meaningful practical application of these models. A null result, however disappointing in some cases, is still potentially valuable and informative.

Marketers, project managers, and business analysts can use this setup to pick out products without peer influence. In any case where homophily is likely to lead to positive correlation in behaviors, a failure to reject the null of zero peer effects is a strong statement. Null results in the presence of a bias of magnitude are more reliable null results. With the results we have presented, there is now a dynamic framework to apply as well (if only as a robustness check). There are mitigating factors which might make a null result less reliable too. Heterogeneity, congestion effects, and omitted variables could make a null result more likely. Future work might extend to matching estimators or simulation techniques to handle these challenges. Yet in the simplest of cases homophily will often lead us to find a putative peer effect that isn't really there. In our case, knowing that Yahoo! Go lacked a positive social externality might have helped Yahoo! end support

for the product earlier, or perhaps to focus their efforts on the most engaged users. Engineering and exploiting contagion in marketing can be a worthy pursuit, but often an expensive strategy to pursue. These models can help decide if resources are better deployed elsewhere.

## References

Angrist, J. D. 2014. "The perils of peer effects," *Labour Economics* (30), Elsevier B.V., pp. 98–108 (doi: 10.1016/j.labeco.2014.05.008).

Angrist, J. D., and Lang, K. 2004. "American Economic Association Does School Integration Generate Peer Effects ? Evidence from Boston ' s Metco Program Published by : American Economic Association Stable URL : http://www.jstor.org/stable/3592836 Does School Integration Generate Peer Effect," *The American Economic Review* (94:5), pp. 1613–1634.

Aral, S., Muchnik, L. E. V, and Sundararajan, A. 2013. "Engineering social contagions: Optimal network seeding in the presence of homophily," *Network Science* (2013), pp. 125–153 (doi: 10.1017/nws.2013.6).

Aral, S., Muchnik, L., and Sundararajan, A. 2009. "Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks.," *Proceedings of the National Academy of Sciences of the United States of America* (106:51), pp. 21544–21549 (doi: 10.1073/pnas.0908800106).

Aral, S., and Nicolaides, C. 2016. "Is exercise contagious? Peer effects in a global health behavior," *MIT working paper*.

Aral, S., and Walker, D. 2011. "Creating Social Contagion Through Viral Product Design : A Randomized Trial of Peer Influence in Networks," (August 2014).

Aral, S., and Walker, D. 2012. "Identifying Influential and Susceptible Members of Social Networks," *Science* (337:July), pp. 337–341 (doi: 10.1126/science.1215842).

Arellano, M., and Bond, S. 1991. "Some tests of specification for panel data: Monte Carlo evidence and an application to employment equations," *The Review of Economic Studies* (58:2), pp. 277–297.

Athey, S., Eckles, D., and Imbens, G. 2015. "Exact P-values for Network Interference," *arXiv:1506.02084v1* (June) (doi: 10.3386/w21313).

Bakshy, E., Hofman, J., Mason, W., and Watts, D. 2011. "Everyone's an influencer: quantifying influence on twitter," *Proceedings of the fourth ACM international conference on Web search and data mining SE - WSDM '11*, pp. 65–74 (doi: doi: 10.1145/1935826.1935845).

Bakshy, E., Rosenn, I., Marlow, C., and Adamic, L. a. 2012. "The role of social networks in information diffusion," *Www2012*, p. 519 (doi: 10.1145/2187836.2187907).

Bapna, R., and Umyarov, a. 2012. "Do Your Online Friends Make You Pay? A Randomized Field Experiment in an Online Music Social Network," *NBER working paper* (June), pp. 1–47 (doi: 10.1287/mnsc.2014.2081).

Barabási, A.-L., and Albert, R. 1999. "Emergence of scaling in random networks," *Science* (286:5439), p. 11 (doi: 10.1126/science.286.5439.509).

Bass, F. M. 1969. "A new product growth for model consumer durables," *Management Science*, pp. 215–227 (doi: 10.1287/mnsc.15.5.215).

Bollinger, B., and Gillingham, K. 2012. "Peer Effects in the Diffusion of Solar Photovoltaic Panels ←," *Marketing Science* (31:6), pp. 900–912 (available at http://pubsonline.informs.org/doi/pdf/10.1287/mksc.1120.0727).

Bramoullé, Y., Djebbari, H., and Fortin, B. 2009. "Identification of peer effects through social networks," *Journal of Econometrics* (150:1), Elsevier B.V., pp. 41–55 (doi: 10.1016/j.jeconom.2008.12.021).

Card, D., and Giuliano, L. 2012. "Peer Effects and Multiple Equilibria in the Risky Behavior of Friends," *NBER working paper* (95:October), pp. 1130–1149 (doi: 10.3386/w17088).

Christakis, N. a, and Fowler, J. H. 2007. "The spread of obesity in a large social network over 32 years.," *The New England journal of medicine* (357:4), pp. 370–9 (doi: 10.1056/NEJMsa066082).

Domingos, P., and Richardson, M. 2001. "Mining the Network Value of Customers," *Proceedings of the*

*Seventh {ACM} {SIGKDD} International Conference on Knowledge Discovery and Data Mining*, pp. 57–66 (doi: 10.1145/502512.502525).

Dupas, P., Duflo, E., and Kremer, M. 2008. "Peer Effects , Teacher Incentives , and the Impact of Tracking : Evidence from a Randomized Evaluation in Citation Accessed Citable Link Detailed Terms Peer Effects , Teacher Incentives , and the Impact of Tracking : Evidence from a Randomized Evaluation," *NBER working paper* (101:5), pp. 1739–1774.

Eckles, D., and Bakshy, E. 2014. "Bias and High-Dimensional Adjustment in Observational Studies of Peer Effects," *Working Paper.*

Fader, P. S., Hardie, B. G. S., and Lee, K. L. 2005. "RFM and CLV: Using Iso-Value Curves for Customer Base Analysis," *Journal of Marketing Research* (XLII:November), pp. 415–430 (doi: 10.1509/jmkr.2005.42.4.415).

Ghose, a., and Han, S. P. 2011. "An Empirical Analysis of User Content Generation and Usage Behavior on the Mobile Internet," *Management Science* (57:9), pp. 1671–1691 (doi: 10.1287/mnsc.1110.1350).

Goldenberg, A., Zheng, A. X., Fienberg, S. E., and Airoldi, E. M. 2010. "A Survey of Statistical Network Models," *Found. Trends Mach. Learn.* (2:2), pp. 129–233 (doi: 10.1561/2200000005).

Goldenberg, J., Libai, B., and Muller, E. 2010. "The chilling effects of network externalities," *International Journal of Research in Marketing* (27:1), Elsevier B.V., pp. 4–15 (doi: 10.1016/j.ijresmar.2009.06.006).

Granovetter, M. 1978. "Threshold Models of Collective Behavior," *American Journal of Sociology* (83:6), pp. 1420–1443 (available at http://www.jstor.org/stable/277811).

Hartmann, W. R., Manchanda, P., Nair, H., Bothner, M., Dodds, P., Godes, D., Hosanagar, K., and Tucker, C. 2008. "Modeling social interactions: Identification, empirical methods and policy implications," *Marketing Letters* (19:3-4), pp. 287–304 (doi: 10.1007/s11002-008-9048-z).

Iyengar, R., Van den Bulte, C., and Valente, T. W. 2011. "Opinion leadership and social contagion in new product diffusion," *Marketing Science* (30:2), pp. 195–212 (doi: 10.1287/mksc.1100.0566).

Jackson, M. O., and Yariv, L. 2007. "Diffusion of Behavior and Equilibrium Properties in Network Games," *The American Economic Review* (97:2), pp. 92–98.

Lee, L. 2007. "Identification and Estimation of Spatial Econometric Models with Group Interactions, Contextual Factors and Fixed Effects," *Journal of Econometrics* (140:2), pp. 333–374 (available at http://www.cemmap.ac.uk/forms/lee_paper.pdf).

Manchanda, P., Xie, Y., and Youn, N. 2008. "The Role of Targeted Communication and Contagion in Product Adoption," *Marketing Science* (27:6), pp. 961–976 (doi: 10.1287/mksc.1070.0354).

Manski, C. F. 1993. "Identification of Endogenous The Reflection Problem," *Review of Economic Studies* (60:60), pp. 531–542 (doi: 10.2307/2298123).

McPherson, M., Smith-lovin, L., and Cook, J. M. 2001. "Birds of a Feather : Homophily in Social Networks," *Annual Review of Sociology* (27), pp. 415–444 (doi: 10.1146/annurev.soc.27.1.415).

Miller, A. R., and Tucker, C. 2009. "Privacy protection and technology diffusion: The case of electronic medical records," *Management Science* (55:7), pp. 1077–1093 (doi: 10.1287/mnsc.1090.1014).

Moffitt, R. 2001. "Policy interventions, low-level equilibria, and social interactions," *Social dynamics* (JUNE 2000), pp. 45–82.

Muchnik, L., Aral, S., and Taylor, S. J. 2013. "Social influence bias: a randomized experiment," *Science* (341:6146), pp. 647–51 (doi: 10.1126/science.1240466).

Page, L., Brin, S., Motwani, R., and Winograd, T. 1998. "The PageRank Citation Ranking: Bringing Order to the Web," *World Wide Web Internet And Web Information Systems* (54:1999-66), pp. 1–17 (doi: 10.1.1.31.1768).

Sacerdote, B. 2001. "Peer effects with random ssignment: Results for Dartmouth Roommates," *Quarterly Journal of Economics* (116:2), pp. 681–704 (doi: 10.1162/00335530151144131).

Schelling, T. 1971. "Dynamic Models of Segregation," *Journal of Mathematical Sociology* (1), pp. 143–186 (available at http://www.stat.berkeley.edu/~aldous/157/Papers/Schelling_Seg_Models.pdf).

Shalizi, C. R., and Thomas, A. C. 2011. "Homophily and Contagion Are Generically Confounded in Observational Social Network Studies.," *Sociological methods & research* (40:2), pp. 211–239 (doi: 10.1177/0049124111404820).

21

Stock, J., and Yogo, M. 2005. "Testing for weak instruments in linear IV regression.," (1; 1:August), pp. 80–108.

Tucker, C. E. 2008. "Identifying formal and informal influence in technology adoption with network externalities," *Management Science* (54:12), pp. 2024–2038 (doi: 10.1287/mnsc.1080.0897).

Watts, D. J. 2002. "A simple model of global cascades on random networks.," *Proceedings of the National Academy of Sciences USA* (99:9), pp. 5766–71 (doi: 10.1073/pnas.082090499).

## Appendix A – Calculation Strategies for $\varrho$

The simplest way to start is to choose either the new network or the old network as the relevant period-level isomorphic graph. We might also define a scalar to indicate the relative "socialness" of each period as well. One way to determine a weighting value between 0 and 1 is to define a row vector of weights for each individual, where the $i^{th}$ element of the row vector is the proportion of the total count of connections belonging to the new network for individual i. This is to say that:

$$\rho_{i,t} = \frac{n_{i,t+1}}{n_{i,t}+n_{i,t+1}}; \; n_{i,t} \text{ is the connection count of person i in time t}$$

We could treat $\varrho$ as another parameter to optimize in the estimation procedure, choosing its value in each period (as a scalar) to minimize the sum of squared residuals. This approach is not discussed in this paper. Lastly we might abandon the convex combination idea entirely, focusing instead on finding the graph union or graph intersection of $F_t$ and $F_{t+1}$. The graph union has the advantage of including all connections as important in peer effects. The graph intersection represents the part of the graph that is stable over time. Most network software packages contain algorithms for calculating these matrices. A combination of both approaches may be worthwhile, taking the "true" network to be some convex combination of the intersection and union graphs.

## Appendix B – First Stage Estimates from Tables 3A and 3B

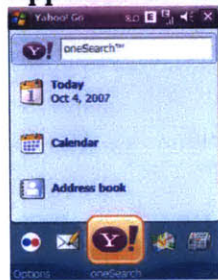| Table 4A – First Stage Dep. Var.: Friend Use | (1) First Stage - Pooled with Lags | (2) First Stage - FE with Lags | (3) First Stage - Pooled with Demographics | (4) First Stage - FE with Demographics |
|---|---|---|---|---|
| Lag Friend-of-Friend Use (1 Day) | 0.0476*** | 0.0306*** | 0.0451*** | 0.0282*** |
| | (0.00922) | (0.00932) | (0.00923) | (0.00933) |
| Lag Friend-of-Friend Use (2 Day) | 0.00325 | 0.0428*** | 0.000941 | 0.0405*** |
| | (0.00920) | (0.00930) | (0.00922) | (0.00932) |
| Friend-of-Friend Age | | | -0.0560 | -0.0644 |
| | | | (0.0496) | (0.0492) |
| Friend-of-Friend Gender | | | -4.947*** | -4.776*** |
| | | | (0.767) | (0.760) |
| Constant | 8.045*** | 9.497*** | 6.908*** | 8.544*** |
| | (0.193) | (0.524) | (1.330) | (1.407) |
| Observations | 57,278 | 57,278 | 57,278 | 57,278 |
| R-squared | 0.199 | 0.215 | 0.201 | 0.216 |

Standard errors in parentheses
*** $p<0.01$, ** $p<0.05$, * $p<0.1$

| Table 4B First Stage<br>Dep. Var.:<br>Friend Use | (1)<br>First Stage - Pooled Full Model | (2)<br>First Stage - FE Full Model |
|---|---|---|
| Lag Friend-of-Friend Use<br>(1 Day) | 0.0327*** | 0.0156* |
| | (0.00926) | (0.00937) |
| Lag Friend-of-Friend Use<br>(2 Day) | -0.0122 | 0.0272*** |
| | (0.00925) | (0.00936) |
| Friend-of-Friend Age | -0.0722 | -0.0803 |
| | (0.0499) | (0.0495) |
| Friend-of-Friend Gender | -3.182*** | -3.114*** |
| | (0.790) | (0.782) |
| Friend-of-Friend Degree | 0.0796 | 0.0543 |
| | (0.158) | (0.156) |
| Friend-of-Friend PageRank | 618.1 | 2,232 |
| | (29,045) | (28,776) |
| Constant | 2.582* | 4.339*** |
| | (1.435) | (1.506) |
| Observations | 57,278 | 57,278 |
| R-squared | 0.205 | 0.220 |

Standard errors in parentheses

*** p<0.01, ** p<0.05, * p<0.1

## Appendix C – Yahoo! Go Screen



Source: Wikipedia