
Capturing Unobserved Correlated Effects in Diffusion in Large Virtual Networks

1 **Elenna R. Dugundji, Michiel van Meeteren**
2 Universiteit van Amsterdam
3 Amsterdam, Netherlands
4 *e.r.dugundji, michielvanmeeteren@gmail.com*

Ate Poorthuis
University of Kentucky
Lexington, KY, USA
atepoorthuis@gmail.com

5 **Abstract**

6 Social networks and social capital are generally considered to be important
7 variables in explaining the diffusion of behavior. However, it is contested
8 whether the actual social connections, cultural discourse, or individual
9 preferences determine this diffusion. Using discrete choice analysis applied
10 to longitudinal Twitter data, we are able to distinguish between social
11 network influence on one hand and cultural discourse and individual
12 preferences on the other hand. In addition, we present a method using freely
13 available software to estimate the size of the error due to unobserved
14 correlated effects. We show that even in a seemingly saturated model, the
15 log likelihood can increase dramatically by accounting for unobserved
16 correlated effects. Furthermore the estimated coefficients in an uncorrected
17 model can be significantly biased beyond standard error margins.

18 19 **1 Introduction**

20 With the onset of ubiquitous social media technology, people leave numerous traces of their
21 social behavior in – often publicly available – data sets. In this paper we look at a virtual
22 community of independent (“Indie”) software developers for the Macintosh and iPhone that
23 use the social networking site Twitter. Using Twitter’s API, we collect longitudinal data on
24 network connections among the Indie developers and their friends and followers
25 (approximately 15,000 nodes) and their use of Twitter client software over a period of five
26 weeks (more than 600,000 “tweets”). We use this dynamic data on the network and user
27 behavior to analyze the diffusion of Twitter client software.

28 Within the Indie community, four prominent software developers have developed Twitter
29 clients (Tweetie, Twittrific, Twittelator, Birdfeed) that compete for adoption within the
30 community. Apart from these Indie Twitter clients, members of the virtual community can
31 choose from a range of clients that are developed outside of the Indie community (for
32 example, Tweetdeck, Twitterfon) as well as the standard Web interface provided by Twitter.
33 Previous qualitative ethnographic evidence for our case study indicates that social networks
34 and social capital are considered to be important factors in explaining the adoption and
35 diffusion of behavior. [1] Using discrete choice analysis applied to longitudinal panel data,
36 we are able to quantitatively test for the relative importance of global cultural discourse,
37 taste-maker influence and other contextual effects, node level behavioral characteristics,
38 socio-centric network measures and ego-centric network measures, individual preferences
39 and social network contagion, in users’ decisions of what client software they choose to
40 interface to Twitter.

41 Importantly, we furthermore demonstrate a method using readily available software to
 42 estimate the size of the error due to unobserved correlated effects in users' choices. This is
 43 critical to test for in any application of multinomial logistic regression where social
 44 influence variables and/or other network measures are used as explanatory variables, since
 45 their use poses a classic case of endogeneity. We show that even in a seemingly saturated
 46 model, the log likelihood of the model fit can increase significantly by accounting for
 47 unobserved correlated effects. Furthermore the estimated coefficients in the uncorrected
 48 model can be significantly biased beyond standard error margins. Failing to account for
 49 correlated effects can yield misleading market share predictions for users' preferences for
 50 Twitter clients.

51 The paper is organized as follows. First a brief review of literature is presented describing
 52 what the paper brings to an existing stream of behavioral modeling research. Next the
 53 understanding of the context of the case study and the insights from the available data lead
 54 us to define nine sets of different kinds of social and individual explanatory variables to
 55 explore in our model, with different functional forms. Estimation results are summarized.
 56 Finally, directions for future research efforts are outlined.

57

58 **2 Discrete choice with social interactions**

59

60 **2.1 Multinomial logit model**

61 Discrete choice analysis allows prediction based on computed individual choice probabilities for
 62 heterogeneous agents' evaluation of alternatives. In accordance with notation and convention in
 63 Ben-Akiva and Lerman [2], the multinomial logit model is specified as follows. Assume a sample
 64 of N decision-making entities indexed $(1, \dots, n, \dots, N)$ each faced with a choice among J_n alternatives
 65 indexed $(1, \dots, j, \dots, J_n)$ in subset C_n of some universal choice set C .

66 The choice alternatives are assumed to be mutually exclusive (a choice for one alternative
 67 excludes the simultaneous choice for another alternative, that is, an agent cannot choose two
 68 alternatives at the same moment in time) and collectively exhaustive within C_n (an agent must
 69 make a choice for one of the options in the agent's choice set). In general the composite choice set
 70 C_n will vary in size and content across agents: not all elemental alternatives in the universal choice
 71 set may be available to all agents. For simplicity in this paper however, we will assume that the
 72 choices are available to all agents.

73 Let $U_{in} = V_{in} + \varepsilon_{in}$ be the utility that a given decision-making entity n is presumed to associate with
 74 a particular alternative i in its choice set C_n , where V_{in} is the deterministic (to the modeler) or so-
 75 called "systematic" utility and ε_{in} is an error term. Then, under the assumption of independent and
 76 identically Gumbel distributed disturbances ε_{in} , the probability that the individual decision-making
 77 entity n chooses alternative i within the choice set C_n is given by:

$$\begin{aligned}
 P_{in} &\equiv P_n(i | C_n) = \Pr(V_{in} + \varepsilon_{in} \geq V_{jn} + \varepsilon_{jn}, \forall j \in C_n) \\
 &= \Pr\left[V_{in} + \varepsilon_{in} \geq \max_{j \in C_n} (V_{jn} + \varepsilon_{jn})\right] = \frac{e^{\mu V_{in}}}{\sum_{\forall j \in C_n} e^{\mu V_{jn}}}
 \end{aligned}$$

79 where μ is a strictly positive scale parameter which is typically normalized to 1 in the multinomial
 80 logit model.

81 The systematic utility is commonly assumed to be defined by a linear-in-parameters function of
 82 observable characteristics \mathbf{S}_n of the decision-making entity and observable attributes \mathbf{z}_{in} of the
 83 choice alternative for a given decision-making entity:

$$84 \quad V_{in} = h_i + V(\mathbf{S}_n, \mathbf{z}_{in}) = h_i + \gamma_i' \mathbf{S}_n + \zeta_i' \mathbf{z}_{in}$$

85 The term h_i is a so-called "alternative specific constant" (ASC), as good practice to explicitly
 86 account for any underlying bias for one alternative over another alternative. In other words, h_i
 87 reflects the mean of $\varepsilon_{jn} - \varepsilon_{in}$, that is, the difference in the utility of alternative i from that of j when

88 all else is equal. Since it is the difference that is relevant, for a general multinomial case with J
89 alternatives we can define a set of at most $J - 1$ alternative specific constants.

90 The terms $\gamma_i = [\gamma_{i1}, \gamma_{i2}, \dots]'$ and $\zeta_i = [\zeta_{i1}, \zeta_{i2}, \dots]'$ are vectors of unknown utility parameters
91 respectively corresponding to the relevant observable agent characteristics S_n , and observable
92 agent-specific attributes z_{in} of the choice alternatives. In general the utility parameters may take
93 alternative specific values, however when there is no variation of the agent characteristics S_n
94 across the choice alternatives, we can define a set of at most $J - 1$ vectors of alternative specific
95 coefficients for the case of the γ_i .

96

97 **2.2 Social interactions**

98 An outstanding challenge in discrete choice analysis is the treatment of the interdependence
99 of various decision-makers' choices [3,4]. Brock and Durlauf [5] introduce social
100 interactions in multinomial discrete choice models by allowing a given agent's choice for a
101 particular alternative to be dependent on the overall share of decision makers who choose
102 that alternative. If the coefficient on this interaction variable is close to zero and not
103 important relative to other contributions to the utility, then the distribution of decision-
104 makers' choices will not effectively change over time in relation to other decision-makers'
105 choices. However, if the coefficient on this interaction variable is positive and dominant
106 enough relative to other contributions to utility, there may arise a runaway situation over
107 time as all decision-makers flock to one particularly attractive choice alternative. In short,
108 the specification captures social feedback between decision-makers that can potentially be
109 reinforcing over the course of time. In diverse literature this is referred to as a social
110 multiplier, a cascade, a bandwagon effect, imitation, contagion, herd behavior, etc. [6]

111 We introduce a social feedback effect among agents by allowing the systematic utility V_{in} to
112 be a linear-in-parameter β first-order function of the proportion x_{in} of a given decision-
113 maker's reference entities who have made this choice. Our model differs from the Brock and
114 Durlauf model in that we consider *non-global* interactions. Agents see different proportions,
115 depending on who their particular reference entities are. Additionally, we also consider
116 various socio-centric and ego-centric network measures and other explanatory variables as
117 contributions to the utility.

118

119 **2.3 Endogeneity**

120 One econometric issue that arises in empirical estimation of social interactions in discrete
121 choice models using standard multinomial logistic regression however, is that the error terms
122 are assumed to be identically and independently distributed across decision-makers. It is not
123 obvious that this is in fact a valid assumption when we are specifically considering
124 interdependence between decision-makers' choices. We might reason that if there is a
125 systematic dependence of each decision-maker's choice on an explanatory variable that
126 captures the aggregate choices of other decision-makers who are in some way related to that
127 decision-maker, then there might be an analogous dependence in the error structure.
128 Otherwise said, the same unobserved effects might be likely to influence the choice made by
129 a given decision-maker as well as the choices made by those in the decision-maker's
130 reference group, which is a classic case of endogeneity. The results and coefficients of such
131 a model are likely to be biased. To try to separate out effects, it is therefore first and
132 foremost critically important to begin with an as well-specified model as possible, making
133 use of relevant available explanatory variables. [7]

134 Dugundji and Walker [8] illustrate issues in the empirical estimation of a discrete choice
135 model with network interdependencies using mixed generalized extreme value model
136 structures with pseudo-panel data. Several modeling strategies are presented to highlight
137 hypothesized interaction effects. In absence of true panel data on interaction between
138 identifiable decision-makers, they use a priori beliefs about the social and spatial dimension
139 of interactions to formulate the connectivity of the network and use socioeconomic data for
140 each respondent as well as the geographic location of each respondent's residence to define
141 aggregate interactions by grouping agents into geographic neighborhoods and into
142 socioeconomic groups where the influence is assumed to be more likely. Technically,

143 however, interactions between identifiable decision-makers may also be modeled using the
 144 approach described given the availability of suitable data.

145 In our empirical case study on adoption of Twitter clients, we do indeed have available data
 146 on which identifiable agents (Twitter users) plausibly influence other identifiable agents'
 147 choices, and furthermore we have longitudinal panel data observing repeated choices by
 148 agents over time. In this paper with such rich data, we continue this exploration of issues in
 149 the empirical estimation of discrete choice models with social interactions. Since our data is
 150 fairly large -more than 10,000 agents- we argue that the effect of unobserved correlated
 151 effects as perceived by any given agent is normally distributed, but is the same for that agent
 152 over the fairly short time period of the data collection. This simplified assumption allows us
 153 to specifically control for correlations in the error structure, through the use of mixed
 154 multinomial logit models with panel effects. [9].

155

156 **2.4 Capturing unobserved correlated effects**

157 Suppose each agent n makes a sequence of choices at a number of points in time indexed
 158 $(1, \dots, t, \dots, T_n)$. For our case study, we will consider a general case where the number T_n of
 159 decision-making moments per agent varies across agents. We introduce an additive, normally-
 160 distributed agent-specific error term for each alternative i as follows:

$$161 \quad U_{int} = V_{int} + \varepsilon_{int} + \sigma_i \xi_{in} ; \quad \xi_n \sim N(0, I)$$

162 Conditional on ξ_n , the probability that agent n makes a particular sequence of choices over time
 163 (i_1, \dots, i_{T_n}) is given by the product of the probabilities for agent n making each individual choice i_i :

$$164 \quad P_n(i_1, \dots, i_{T_n} | \xi_n) = \prod_{\forall t \in T_n} \frac{e^{\mu(V_{int} + \sigma_i \xi_{in})}}{\sum_{\forall j \in C_n} e^{\mu(V_{jnt} + \sigma_j \xi_{jn})}}$$

165 The unconditional user choice probability is the integral of this product over all values of ξ_n

$$166 \quad P_n(i_1, \dots, i_{T_n}) = \int_{\xi_n} \prod_{\forall t \in T_n} \frac{e^{\mu(V_{int} + \sigma_i \xi_{in})}}{\sum_{\forall j \in C_n} e^{\mu(V_{jnt} + \sigma_j \xi_{jn})}} N(0, I) d\xi_n$$

167

168 **2.5 Econometric estimation with simulation**

169 The unconditional choice probability is approximated through simulation for any given value of
 170 ξ_{in} as follows:

- 171 1) Draw a vector of values of ξ_n from $N(0, I)$ for each alternative in the choice set C_n , and label
 172 this ξ_n^r with the superscript $r = 1$ referring to the first draw
- 173 2) Calculate the conditional user choice probability for the particular sequence of choices made by
 174 agent n with this draw
- 175 3) Repeat steps 1 & 2 for R total number of draws and average the results

$$176 \quad \hat{P}_n(i_1, \dots, i_{T_n}) = \frac{1}{R} \sum_{r=1}^R \prod_{\forall t \in T_n} \frac{e^{\mu(V_{int} + \sigma_i \xi_{in}^r)}}{\sum_{\forall j \in C_n} e^{\mu(V_{jnt} + \sigma_j \xi_{jn}^r)}}$$

177 If the estimated coefficients σ_i can be shown to be statistically insignificant, we assume that the
 178 hypothesized endogeneity has negligible effect.

179

180 **3 Modeling the effects**

181 This paper studies the diffusion of Twitter clients within the Indie community. Based on
182 earlier research [10] we were able to determine a community of Indie developers that are
183 actively using Twitter, using a mixed method community detection approach. For this
184 community we use Twitter's publicly available API to gather data on network connections
185 and actual messages sent. For 39 days, from 9 August until 16 September 2009, we harvested
186 tweets and network connections on a daily basis for each of the nodes in the community.

187 Based on a review of the case study and the data [11], we expect client choice to be
188 influenced by a number of distinct dimensions. Generally, social networks or social capital
189 are considered to be important variables in explaining the adoption and diffusion of
190 behavior. However, it is debatable to what extent the actual social connections, the global
191 cultural discourse, and individual preferences influence this adoption and diffusion. Through
192 our modeling of the effects, we can try to test the different hypotheses.

193 We distinguish between four Indie clients (Tweetie, Twitterrific, Birdfeed and Twittelator),
194 two popular non-Indie clients (Twitterfon and Tweetdeck) and the default Twitter web
195 interface ("Web"). In addition, we employ a choice alternative, "Other" that serves as a
196 baseline reference for the modeling. The "Other" category is highly heterogeneous and
197 consists of more than 3500 clients that have relatively small market share ($< 1\%$). On the
198 basis of data we proceed to construct the following nine sets of different kinds of social and
199 individual explanatory variables to explore in our model.

200

201 **3.1 Contextual effects: taste maker influence**

202 We start with exploring the contextual effect of whether or not a user in the community is
203 connected to professional independent tech blogger John Gruber. Since Gruber promotes
204 different clients to different extents [11], we are interested to see if the clients he promotes
205 most favorably are used more often by the users connected to him. We operationalize this
206 dummy variable in two different ways: if a user "follows" Gruber (ie. user receives tweets
207 from Gruber); and if there is a reciprocal link with Gruber.

208

209 **3.2 Contextual effects: developer influence**

210 Next, we are interested in the contextual effect of whether or not a user in the community is
211 connected to a Twitter client developer [11] as follows: Clients developed by "Indies":
212 Tweetie (Loren Brichter, Atebits); Twitterrific (Craig Hockenberry, Iconfactory); Twittelator
213 (Andrew Stone, Stone Design); Birdfeed (Buzz Andersen, SciFi HiFi); Clients developed by
214 others: TweetDeck (Iain Dodsworth, TweetDeck); TwitterFon (Kazuho Okui, Naan Studio).
215 We operationalize each of these dummy variables in two different ways: if a user "follows"
216 the developer (ie. user receives tweets from the developer); and if there is a reciprocal link
217 with the developer (ie. the link with the developer is especially strong).

218

219 **3.3 Behavioral characteristics: power users**

220 Since the Twitter clients have very different features, we might expect users who tweet a lot
221 to prefer different kinds of clients than users who tweet less frequently. We operationalize
222 this variable in four different ways: number of tweets sent by a user during observation
223 period; "status count" (total tweets sent by a user during their entire history); number of
224 tweets sent by a user prior to observation period (ie. giving emphasis of how active the user
225 was in the past and how long the user has been using Twitter); and finally, the ratio of tweets
226 sent by a user during observation period to total tweets sent during their entire history.

227

228 **3.4 Network measures: central users**

229 As per our review of the importance of social media networks for "echo-chamber"
230 marketing, we are interested in whether a user's position in the community affects their
231 client choice. We compute five classic network centrality measures: in-degree centrality (the
232 number of a user's "friends" in sample, ie. from whom tweets are received); out-degree

233 centrality (the number of a user’s “followers” in sample, ie. to whom tweets are sent);
234 closeness centrality (sum of distances from a user to all other users, giving an indication of
235 the expected time until arrival for information that might be flowing through the network);
236 betweenness centrality (how often a user lies along the shortest path between two other
237 users, giving an indication of access to diversity of information); and finally, eigenvector
238 centrality (measures if a user is connected to many users who are themselves well connected,
239 identifying users in centers of cliques).

240

241 **3.5 Network measures: extended user in-degree**

242 In order to test the relative importance of the exposure to information flowing through the
243 wider Twitter universe outside of the Indie community, we explore three extra network
244 measure variables: the total number of a user’s “friends” in the entire Twitter universe, ie.
245 from whom a given user in principle receives tweets; the number of users outside the
246 community from whom a given user in principle receives tweets; and finally, the ratio of
247 users inside sample from whom a given user receives tweets to their total “friends” in the
248 Twitter universe.

249

250 **3.6 Network measures: extended user out-degree**

251 Similarly, in order to test the relative opportunity to influence other users in the wider
252 Twitter universe outside of the Indie community, we explore three extra network measure
253 variables: the total number of a user’s “followers” in the entire Twitter universe, ie. to whom
254 a given user in principle sends tweets; the number of users outside the community to whom a
255 given user in principle sends tweets; and finally, the ratio of users inside sample to whom a
256 given user sends tweets to their total “followers” in the Twitter universe.

257

258 **3.7 Temporal effects: individual preferences**

259 We operationalize individual preference by constructing an alternative-specific relative
260 individual cumulative lag variable. For each tweet, we count how often the sending user has
261 been using each client in the seven days prior to sending the tweet resulting in an absolute
262 cumulative lag variable. For each client, we then convert this absolute frequency to a relative
263 cumulative lag variable indicating that client's use relative to how often that user has been
264 using other Twitter clients in the past seven days. This individual preference variable shows
265 how “sticky” a particular client has been for a user in the past seven days. This individual
266 past behavior is likely to be a predictor of client choice for the next tweet, capturing
267 complex UI preferences which we as researchers were not able to measure directly.

268

269 **3.8 Temporal effects: social network contagion**

270 To operationalize network influence we use the absolute cumulative lag variable as a basis.
271 For each tweet, we count how often all users that the sender of that specific tweet is
272 following use each client in the seven days prior to sending that the tweet. We convert the
273 absolute frequency to an alternative specific relative network influence variable that
274 indicates how often each client has been used relative to all other clients by all users that the
275 sender of the tweet is following (ie. receiving information from). This can entail specific
276 mentions of a client in a tweet but also more implicit or tacit knowledge about which client
277 is popular or deemed useful within that user's social network. We argue that this usage by
278 “friends” might influence client choice by either specific mentions of a client in Tweets or
279 by the effect of tacit knowledge encoded within a user's social network.

280

281 **3.9 Global influence**

282 The cultural discourse on what is popular within the entire Indie community is
283 operationalized by a set of alternative specific constants (ASC). Amongst things such as
284 price and the impact of media exposure, we argue that this effectively captures global
285 influence. It indicates the popularity of an alternative relative to all other alternatives during
286 the entire sample period, after controlling for all other effects.

287 **4 Results**

288 All models are estimated using the freely available optimization toolkit Biogeme
 289 (<http://biogeme.epfl.ch>) developed by Bierlaire. We begin by estimating a baseline
 290 multinomial logit model with alternative specific constants only, representing global bias.
 291 The log likelihood, number of estimated parameters and adjusted rho-squared are given in
 292 the first line of Table 1.

293

294 Table 1: Log likelihood tests for incremental model specifications

295

Nr	Log Likelihood	Est. Par	Rho Sq	-2[L _R -L _U]	$\chi^2(0.1)$	p-Value
1	-968350.6	7	0.262	--	--	--
2	-954368.7	14	0.272	27964	18.5	0.000
3	-568721.3	21	0.566	771295	18.5	0.000
4	-567945.3	28	0.567	1552	18.5	0.000
5	-566798.7	34	0.568	2293	16.8	0.000
6	-562010.9	41	0.571	9576	18.5	0.000
7	-561154.7	48	0.572	1712	18.5	0.000
8	-560664.6	55	0.572	980	18.5	0.000
9	-559662.6	62	0.573	2004	18.5	0.000
10	-559546.0	69	0.573	233	18.5	0.000
11	-452048.1	76	0.655	214996	18.5	0.000

1: Baseline model with alternative-specific constants only; 2: + Social network contagion (sq root); 3: + Lagged individual preferences (sq root); 4: + Follows Gruber; 5: + Follows developer; 6: + Frequency tweets during observation period (sq root); 7: + Eigenvector centrality (sq root); 8: + Closeness centrality; 9: + Ratio in-degree to total friends in Twitter (sq root); 10: + Ratio out-degree to total followers in Twitter (sq root); 11: + Estimated user-specific error component

296

297 Next we test one-by-one each of the explanatory variables defined in Section 3.1-3.8. In
 298 cases where the variables are continuous (ie. for all cases except for the dummy variables in
 299 section 3.1 and 3.2), we also test linear, quadratic and square root forms of these variables.
 300 Based on log likelihood tests compared to the baseline model and t-tests on the estimated
 301 coefficients [2], we identify the best fitting variables per category. For example, the
 302 dummies defined as “follows Gruber” and “follows developer” are more significant than
 303 their respective forms “reciprocal link with Gruber” and “reciprocal link with developer”;
 304 the most significant centrality measures are closeness and square root of eigenvector
 305 centrality, etc. The interested reader is referred to [11] for details and interpretation.

306 Having determined the best fitting variables and their respective functional forms, we then
 307 add the variables incrementally to the model, testing the improvement in log likelihood at
 308 each step. This is important to do, since variables that may have been significant when
 309 included in the model specification on their own, might no longer be significant when
 310 included together due to significance being shared between variables. The results are
 311 reported in lines 2-10 of Table 1. Each successive specification adds seven new parameters
 312 to the model (with the exception of “follows developer” where there are six since the Web
 313 alternative does not have a third party developer), as our data is rich and extensive enough to
 314 support alternative-specific definitions of the variables. In our case study, each new set of
 315 variables significantly improves the log likelihood (p-value of 0.000).

316 Finally, we include the normally-distributed user-specific error terms as in Section 2.4. We
 317 test the robustness of results using three different optimization algorithms for the
 318 maximization of the log likelihood, each with ten different random seeds for generating the
 319 draws. We use the estimated coefficients from the model in line 10 of Table 1 as a starting
 320 point for these 30 estimation runs with 50 draws, and then use the results with 50 draws in
 321 turn as the starting point for another 30 estimation runs with 200 draws, etc., for increasing
 322 number of draws, until the results stabilize across the random seeds for the three different

323 optimization algorithms. Accounting for the unobserved correlated effects gave a dramatic
324 jump in log likelihood as seen in line 11 of Table 1. The estimated coefficients in the final
325 model in line 11 were also significantly different beyond standard error margins for 63 of 69
326 variables in the model in line 10. [11] Failing to account for unobserved correlated effects
327 can thus yield misleading market share predictions for users' preferences for Twitter clients.

328

329 **7 Conclusions and recommendations**

330 A prominent approach to studying the dynamics of networks and behavior stems from a
331 growing stream of research on stochastic actor-based models. See Snijders, van de Bunt, and
332 Steglich [12] for a tutorial. With the large data in our case study however, these established
333 methods are not tractable. The alternative approach we discuss in this paper allows us to
334 apply other freely available, open source, existing software for the estimation of the models.
335 In so doing, we hope to stimulate researchers and practitioners to adopt these techniques
336 when using large data sets of more than 1000 nodes due to the relatively lower entry barrier
337 than could be the case if dedicated code would need to be written or if expensive software
338 would need to be purchased. An interesting direction for further discrete choice research on
339 diffusion in large networks may be combining the approach of Aral, Muchnik and
340 Sundararajan [13] for distinguishing causal effects using propensity score matched sample
341 estimation in dynamic networked settings, with the present work accounting for unobserved
342 correlated effects.

343

344 **References**

345 [1] M. van Meeteren (2008). Indie fever: The genesis, culture and economy of a community of
346 independent software developers on the Macintosh OS X platform. A Sofa Publication [On-line].
347 Available: <http://www.madebysofa.com/indiefever>

348 [2] M. Ben-Akiva and S. R. Lerman (1985) *Discrete Choice Analysis: Theory and Application to*
349 *Travel Demand*. Cambridge, MA: MIT Press.

350 [3] D. McFadden (2001) Economic choices. *American Economic Review* 91(3): 351-378.

351 [4] D. McFadden (2010) Sociality, rationality, and the ecology of choice. In S. Hess and A. Daly
352 (eds.), *Choice Modelling: The State-of-the-Art and the State-of-Practice*. Bingley, UK: Emerald Group
353 Publ. Ltd.

354 [5] W. A. Brock and S. N. Durlauf (2002) A multinomial choice model of neighborhood effects.
355 *American Economic Review* 92(2): 298-303.

356 [6] C. F. Manski (1995) *Identification Problems in the Social Sciences*. Cambridge, MA: Harvard
357 Univ. Press.

358 [7] E. R. Dugundji (2012) *Socio-Dynamic Discrete Choice: Theory and Application*. Ph.D. manuscript.
359 Universiteit van Amsterdam, Netherlands.

360 [8] E. R. Dugundji and J. L. Walker (2005) Discrete choice with social and spatial network
361 interdependencies: An empirical example using mixed generalized extreme value models with field
362 and panel effects. *Transportation Research Record* 1921:70-78.

363 [9] K. E. Train (2009) *Discrete Choice Methods with Simulation, 2nd ed.* New York: Cambridge
364 University Press.

365 [10] M. Meeteren, A. Poorthuis, and E. R. Dugundji (2010) Mapping communities in large virtual
366 social networks, IEEE International Workshop on Business Applications of Social Network Analysis,
367 Bangalore, India. [Online]. Available: <http://dx.doi.org/10.1109/BASNA.2010.5730297>.

368 [11] E. R. Dugundji, A. Poorthuis, and M. van Meeteren (2012) Capturing unobserved correlated
369 effects in diffusion in large virtual networks: Distinguishing individual preferences, social connections
370 and cultural discourse influence on the adoption of Twitter clients, unpublished.

371 [12] T. A. B. Snijders, G. G. van de Bunt, and C. E. G. Steglich (2010) Introduction to stochastic
372 actor-based models for network dynamics. *Social Networks* 32(1): 44-60.

373 [13] S. Aral, L. Muchnik and A. Sundararajan (2009) Distinguishing influence-based contagion from
374 homophily-driven diffusion in dynamic networks. *Proc. Nat'l Acad. Sci.* 106(51): 21544-21549.