



United Nations
University

WIDER

World Institute for Development Economics Research

Discussion Paper No. 2002/57

Risk-Sharing and Endogenous Network Formation

Joachim De Weerd*^{*}

June 2002

Abstract

In economic literature insurance networks are often treated as exogenous institutions. Frequently, the assumption is made that some clearly identifiable group (e.g. ‘the whole village’ or ‘the extended family’) constitutes an insurance network. Still, theory suggests that the formation of insurance links depends on a myriad of factors related to smooth information flows, norms, trust, the ability to punish, discount rates, group size and the potential gains of cooperation (e.g. how correlated income streams are and the amount of other income smoothing strategies available). Recent research has shown that risk is more likely to be shared in small, tightly knit networks, which do not necessarily coincide with a collection of households delineated on the basis of just one factor like kinship or place of residence. However, so far only few empirical attempts have been made at investigating the determinants of network formation. This paper suggests an approach by which one might present stylized facts on endogenous network formation. Applying it to a data set collected in a Haya village in rural Tanzania, we find that kinship, geographical proximity, the number of common friends, clan membership, religious affiliation and wealth strongly determine the formation of risk-sharing networks. Our data suggest that poor households have less dense networks than the rich, making them more vulnerable in the face of idiosyncratic risk.

Keywords: risk-sharing, networks, group formation

JEL classification: O12, O17, Z13

Copyright © UNU/WIDER 2002

*Katholieke Universiteit, Leuven

This study has been prepared within the UNU/WIDER project on Insurance Against Poverty, which is directed by Dr Stefan Dercon.

UNU/WIDER gratefully acknowledges the financial contribution to the project by the Ministry for Foreign Affairs of Finland.

Acknowledgements

Research funded by the National Science foundation (FWO-Vlaanderen). I am indebted to Abigail Barr, Stefan Dercon, Micheal Kevane, Markus Goldstein, Els Lievois, Katleen Van den Broeck and participants at the UNU/WIDER workshop on Insurance Against Poverty, the NEUDC 2001 conference at Boston University, the 'understanding poverty and growth in Sub-Saharan Africa' conference at Oxford University and a seminar at the University of Dar es Salaam for their comments and help. I thank Philemon Charles, Augustina John, Obadiah Kyakajumba, Respichius Mitti, George Musikula, Isaya Mukama, Adelina Rwechungula, and Taddeo Rweyemamu for their excellent work on the field, Julius Majula for taking charge of the data entry and all the people of the community of Nyakatoke for their continuous hospitality and willingness to provide us with the data. This paper based on my doctoral thesis on 'Social Networks, Transfers and Insurance in Developing Countries', supervised by Stefan Dercon. For comments please contact joachim_dw@yahoo.com.

UNU World Institute for Development Economics Research (UNU/WIDER) was established by the United Nations University as its first research and training centre and started work in Helsinki, Finland in 1985. The purpose of the Institute is to undertake applied research and policy analysis on structural changes affecting the developing and transitional economies, to provide a forum for the advocacy of policies leading to robust, equitable and environmentally sustainable growth, and to promote capacity strengthening and training in the field of economic and social policy making. Its work is carried out by staff researchers and visiting scholars in Helsinki and through networks of collaborating scholars and institutions around the world.

UNU World Institute for Development Economics Research (UNU/WIDER)
Katajanokanlaituri 6 B, 00160 Helsinki, Finland

Camera-ready typescript prepared by Jaana Kallioinen at UNU/WIDER
Printed at UNU/WIDER, Helsinki

The views expressed in this publication are those of the author(s). Publication does not imply endorsement by the Institute or the United Nations University, nor by the programme/project sponsors, of any of the views expressed.

ISSN 1609-5774
ISBN 92-9190-238-1 (printed publication)
ISBN 92-9190-239-X (internet publication)

1 Introduction

In much of the economic literature an insurance network is treated as an exogenous institution. Frequently, the assumption is made that all households who are member of an easily identifiable group (e.g. the village, the extended family, etc.) form one single network. There are, however, both theoretical and empirical grounds for being sceptical about this approach. Bala and Goyal (2000) theoretically model network formation as a non-cooperative game. Although the specifics of their model are geared towards explaining the formation of information networks, the basic principle can also be applied to risk-sharing networks: self-interested individuals can form or destroy links with others, trading off the costs and the benefits of doing so. Numerous factors will determine whether agents are able to exploit the gains of cooperation in the context of an informal insurance market. In an economy with heterogeneous agents these factors will differ across dyads.¹ This makes some pairs of households more suitable insurance partners than others.

First, there are factors related to information. Agents who have smooth information flows between each other are more likely to enter into an insurance arrangement. This means, for example, that we expect close neighbours or households engaged in similar income generating activities to form insurance links with each other. Note that even between members of the same village or extended family information flows are likely to be heterogeneous.

Second, there are factors related to trust, norms and the ability to punish (Platteau 1991 and Fafchamps 1992). Kinship, clan membership and religious affiliation might be important in this respect, because they help to impose strict norms on members. Deviant behaviour of group members can be punished with disgrace or ostracism. *Ceteris paribus*, this creates an incentive to form links within one's group.

Third, if dyads are heterogeneous with respect to the correlation of their income flows, then the potential gains of cooperation may differ greatly across dyads. Two households engaged in different activities may have weakly correlated income streams and may thus be better insurance partners (if we abstract from any informational concerns). Households, who have similar activities or, in an agricultural setting, are close neighbours, might have more covariate income streams and are less suitable as insurance partners. Grimard (1997) points to the trade-off that agricultural households are forced to make regarding the geographical proximity of their insurance partners. Informational flows are smooth between close neighbours, but income streams are likely to be covariate. Households living further away might have less covariate incomes, but informational problems can be large. Our data suggest that the information constraints outweigh the non-covariance of income.

Fourth, non-cooperative, game theoretic models (e.g. Coate and Ravallion 1993, Spinnewyn and Wiyaya 1996) stress the importance of the discount factor. Agents need to be sufficiently concerned about future income flows to engage in risk-sharing arrangements.

¹ A dyad is a pair of households. When we say 'across dyads', we mean across all possible combinations of two households in the village.

The fact that the creation of insurance links depends on such a wide range of factors and that dyads are heterogeneous with respect to these factors suggests that insurance networks will not be identical to groups delineated on the basis of just one factor. For example, there is no *a priori* reason to assume that the whole village forms a network. Indeed, in many societies a single village is spread over a substantial area and information flows are, *ceteris paribus*, better between close neighbours than between villagers living, say, 1 km apart. Typically there will also be different clans, religions, kinship and professional groups within a single village, making the likelihood of forging an insurance link unequally distributed across all possible dyads in the village.

Researchers have offered other compelling evidence that there are bounds on the size of an insurance group, even when agents are homogenous. Murgai et al. (2000) argue that there are increasing costs to group size. As the network becomes larger the task of coordinating transfers, gathering information and enforcing contracts becomes more difficult. In such an environment, full insurance at village level becomes an extreme case. They back this argument up with an empirical study of water exchanges along irrigation canals in Pakistan. Genicot and Ray (2000) show that one does not even have to impose increasing costs to have bounded group size. They consider a non-cooperative risk-sharing model, which is robust not only to single-person deviations, but also to subgroup deviations. They show that introducing this (quite natural) assumption is sufficient to put bounds on the size of the network. In both these studies, community level (or family level) insurance can be seen as an extreme case.

The little empirical work that has been done on these issues all suggests that networks should not be treated as exogenous. Fafchamps and Lund (2000) and Murgai et al. (2000) find that risk is shared within the confines of small clusters of households and not at the level of the community (and neither at any other clearly demarcated level). Murgai et al. (2000) find that kinship, geographical proximity and the degree of risk exposure are important in explaining network formation. Goldstein (2000) notes that networks are formed among kin, neighbours and gender groups.

This paper builds on these theoretical and empirical findings and suggests an approach by which one can present stylized facts on the determinants of network formation. Our unit of analysis will be the dyad, a pair of households. Dyads are often studied in the literature on network analysis, because of their intuitive appeal and their versatility in statistical tests (e.g. Wasserman and Faust 1995). Applied to a data set collected in a small Haya village in rural Tanzania, we find that kinship, geographical proximity, the number of common friends, clan membership, religious affiliation and wealth strongly determine the formation of risk-sharing networks.

2 The dyadic data

The basic procedure of the dyadic approach is to make a data set of all possible unique combinations of two households in the village. For each of these pairs we construct a value for the ‘degree of connection’, or ‘the strength of the insurance link’ between them. This will be the endogenous variable, which can be explained by exogenous variables (e.g. the strength of the kinship tie) through multiple regression analysis. This paper does not analyse any issues related to the *direction* of the link.

We make use of data collected in the community of Nyakatoke, a small Haya village in the Kagera region of Tanzania. We interviewed *all* of the 119 households in the community in household interviews and *all* of the 220 adults in these households in individual interviews. Although some of the data has been collected at individual level, our analysis will be done at household level. Note that in a village of 119 households there are $\binom{119}{2} = 7021$ unique combinations of two households (also called dyads) possible.

We want to attach a value of degree of connection to each of these 7021 dyads. Because of time constraints, we did not query each of the 220 individuals about their relation with each of the other 119 individuals in the survey. Instead we asked them ‘*Can you give a list of people from inside or outside of Nyakatoke, who you can personally rely on for help and/or that can rely on you for help in cash, kind or labour*’.

The persons mentioned are called the ‘network members’ or the ‘network partners’. The respondents listed a total of 1126 network members. About two thirds of them (738 in total) live in Nyakatoke. The other one third (388 in total) live outside of Nyakatoke. Given the set-up of our approach, we are only able to make use of the 738 links inside Nyakatoke.

We want to make two remarks concerning the weight that will be attached to each link. First, although our analysis is at household level,² one might argue that if more than one member of the household mentions a particular other household, this link should have a higher weight. Second, respondents often unilaterally mention each other as network partners.³ This should not be taken to mean that there are false expectations or the relation is not reciprocal. Remember that the question was framed as ‘who can you rely on and/or who relies on you’, so no directional meaning can be attached to it. In this paper, we take the view that the interviews administered to each side of the dyad complement one another. We consider unilaterally mentioned links to be weaker than bilaterally mentioned links.

Taking account of the two previous remarks, we construct three different specifications for the LHS variable, summarized in the first three rows of Table 1. First, we have the total number of links that were reported between two households. For example, if two members of household A mention someone from household B and one member of B mentions a member of A, then the total number of links is three. The second measure also counts the links, but now under the assumption that each household can only send one link to (and receive one link from) each of the other 118 households. If several members of household A send a link to household B, then this is counted as a single link. In the above example this variable would be equal to two. The variable is equal to two when the link between A and B is reciprocal, equal to one when it is unilateral and equal to zero when neither of the two households mentions each other. The third specification is simply a dummy indicating whether there exists at least one link between the two households. We will perform an ordinal logit regression for the first two variables and a logit regressions for the third.

² An analysis at individual level would have 24090 dyads, of which 23352 have the endogenous variable equal to zero. This would give us too little variation for the econometric analysis.

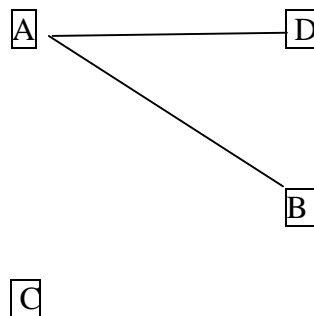
³ Fafchamps and Lund (2000) find the same for households in the Philippines.

Table 1
Four different specification of the endogenous variable

Endogenous variable	Definition	Distribution across the 7021 dyads	
Total number of links reported by individuals	The number of times a member of household A mentions a member of household B + the number of times a member of household B mentions a member of household A	0	6531
		1	308
		2	129
		3	42
		4	9
		5	2
Number of household level links	0 if there is no link between the households, 1 if there is a unilateral link and 2 if there is a reciprocal link.	0	6531
		1	350
		2	140
Dummy	0 if there is no link between the households 1 if there is at least one link	0	6531
		1	490
Geodesic distance	Shortest distance between the households on the network graph	1	490
		2	1996
		3	2900
		4	1275
		5	360

Source: Nyakatoke Household Survey

The fourth measure for our LHS variable is the ‘geodesic distance’ between two households. Because we interviewed every household in the village we can, quite literally, draw the complete network inside the village, with lines joining the households which have mentioned each other as network members (this graph is given in Appendix 1). The ‘geodesic distance’ between two households is the minimum amount of steps one has to take to go from one household to the other on the network graph (Wasserman and Faust 1995).



For example, in the picture given above A has a link with B and D; C is not linked to anyone. The geodesic distance is 1 between A and B and between A and D. It is two between B and D and it is infinity between C and any other household. The geodesic distance is more convincing for those who believe that one does not only benefit from one’s own direct network members, but also from the vast network lying behind one’s own direct network members. Of course, there are costs attached to each node that is crossed in the network (Bala and Goyal 2000). The geodesic distance is an ordinal

measure for these costs and thus a natural measure of how well connected two households are. The last row of Table 1 shows how the geodesic distance is distributed across the dyads. We use an ordinal logit to explain the geodesic distance between two households.

The survey has data on religion, clan, schooling, income generating activities, age of the household head and wealth. Table 2 and 3 present summary statistics for these variables. Out of the three religions in the village, the Muslims form the smallest group. There are 26 clans in Nyakatoke. The Bayango form the biggest clan (23 members), followed by the Basimba (20 members), the Bahimba (12 members) and the Bahunga (10 members). There are 10 clans which have only one household in Nyakatoke. Schooling is very low in this village, with only 72 household heads having completed primary school. There are 93 households in which at least one member has completed primary education. Almost all households are engaged in crop production. Note that very few people get any income out of assets. Livestock and land holdings are good indicators of wealth in Haya society (Reining 1967). Table 3 shows how unequally distributed wealth is in the village, especially livestock.

Table 2
Distribution of attribute variables

Variable	Category	Frequency	<i>N</i>
Distribution of religion	Muslim	24	119
	Lutheran	46	
	Catholic	49	
The number of clans with the specified number of members*	1 member	11	119
	2 members	5	
	3 members	2	
	4 to 10 members	5	
	12 to 23 members	3	
Household has at least one member who has completed primary education	No	26	119
	Yes	93	
Income generating activities. Number of households engaged in each activity	casual labour	57	116
	trade	41	
	crops	108	
	livestock	31	
	assets	8	
	processing	45	
	other off-farm	40	

* The exact distribution, specified per clan, can be found in Mitti and Rweyemamu (2000).

Source: Nyakatoke Household Survey.

Table 3
Distribution of attribute variables, quintiles

Quintile	Mean age of household head (years)	Mean livestock value (TSh)	Mean area of land (ha)
1 st (lowest)	27	0	0.29
2 nd	35	1867	0.62
3 rd	41	7400	0.98
4 th	54	23544	1.45
5 th highest	70	254018	3.45
Total	45	53354	1.34
<i>N</i>	119	119	119

Source: Nyakatoke Household Survey.

Table 4
Distribution of relational variables across the dyads

Variable	Category	Frequency	<i>N</i>
Kinship	Parents, children and siblings	109	7021
	Nephews, nieces, uncles, aunts, cousins, grandparents and grandchildren	102	
	other blood bond	172	
	no blood bond	6638	
Number of common friends	0	4696	7021
	1	1421	
	2	563	
	3	209	
	4	69	
	between 5 and 11	63	
Same religion	Yes	2487	7021
	No	4534	
Same clan	Yes	659	7021
	No	6362	
Do both households have at least one member who has completed primary education?	no, neither of them has	325	7021
	no, only one of them has	2418	
	yes, both of them have	4278	

Source: Nyakatoke Household Survey.

Table 5
Distribution of relational variables across dyads, quintiles

Quintile	Distance (metres)	Activity overlap (%)	Age of household head (years)		Livestock value (TSh)		Land (ha)		No. of respondents
			diff.	max.	diff.	max.	diff.	max.	
1 st (lowest)	161	0.19	3	35	678	927	0.12	0.64	2
2 nd	325	0.45	9	44	6496	9464	0.41	1.05	2
3 rd	466	0.57	16	53	16805	23482	0.75	1.46	2
4 th	645	0.66	24	63	36947	44680	1.40	2.30	2
5 th highest	1016	0.85	40	78	403638	436736	3.39	4.41	3
TOTAL	522	0.53	18	54	92727	99718	1.21	1.94	4
N	6670	6670	7021	7021	7021	7021	7021	7021	7021

Source: Nyakatoke Household Survey.

The above mentioned variables are all *attribute variables*. They can only take on a value for a single household. We cannot sensibly talk about the religion of two households. However, the observations for the regression analysis are dyads, so we can only enter *relational variables* as explanatory variables. There are some variables that are intrinsically relational (e.g. kinship ties) and others that can be created by transforming attribute variables. We enter the following relational variables in our regressions (summary statistics are given in Tables 4 and 5).

2.1 Kinship

Kinship might be important because of history, norms and trust. Typically, one has had a long lasting relationship with family members and, as Hayami and Kikuchi cited in Platteau (1991) put it, ‘performances in past transactions comprise a reliable data set for prediction of future performances’. Furthermore, the family as a group is likely to criticise or punish uncooperative behaviour, thus inducing norms and trust. From a Darwinist point of view, helping family members is good for the expansion of the gene pool. There are data on the kinship ties between all the households in Nyakatoke. We use these to create three dummy variables. The first is equal to one for parents, children and siblings. The second is equal to one for nephews, nieces, uncles, aunts, cousins, grandparents and grandchildren. The third is equal to one for any other blood relation. The reference category is having no blood bond. There are 383 blood bond relations between the households of Nyakatoke. We do not include relations which are through clan membership only. The influence of clan membership in network formation will be captured by a separate variable. The Haya have a patrilineal society, so it would make sense to create kinship dummies that take this into account. However, because the wife usually moves to the husband’s village when she marries, there are very few matrilineal links in the village. For example, in the second category (uncles, cousins etc.) there are only 8 matrilineal links, which is too little to enter as a separate variable.

2.2 Geographical distance

Neighbours have clear information advantages, which we expect to have a positive effect on the emergence of risk-sharing arrangements. On the other hand they are likely to have covariate agricultural risk, which we expect to hamper the formation of a link. Haya homesteads are typically surrounded by dense fields of banana trees, intercropped with other crops. These form a natural barrier with their closest neighbour, who can only be reached by a short walk through the field or along a path. Houses are never built adjacent to each other. All the homesteads were plotted on an electronic map and the distance between each pair was calculated. The average distance between the households is 522 metres; the maximum distance 1738 metres.

2.3 Number of common friends

Households having many common friends are likely to have more information channels (e.g. through the grapevine) between them than households without common friends. It may also be true that these common friends act as observers of the relation between the two households and will frown upon any deviant behaviour. Falling out with someone might impose the extra punishment of falling out with some common friends, especially if you are the one considered to be behaving in the wrong way. We measure this variable as the number of common network members both households have. If A is linked with B, C, D and E and B is linked with C, E, F, G and H, then A and B have two common 'friends' (C and E). In the data this variable ranges from 0 to 11.

2.4 Same religion

One might argue that religious gatherings ease information transfers, but the control and mediating functions the church or mosque have are probably more important. We expect that, *ceteris paribus*, links form within the same religious groups. 'Same religion' is a dummy, which is one if both households are of the same religion and we expect it to have a positive coefficient.

2.5 Same clan

Although it is in the process of losing its significance, the clan is still an important institution in Haya culture, for example in matters regarding land rights. The clan elders can, in effect, function as a court of law. They could easily reprimand younger clan mates when they think their behaviour is bad for the clan. Everybody wants to avoid falling out with their clan. 'Same clan' is a dummy, which is one if both households are of the same clan. We expect this dummy to have a positive sign.

2.6 Education

Do households link up with others who have the same educational achievement as them, or are insurance partners mixed in terms of education? With respect to this, note that there may be inter-household externalities to education, which make it interesting for a non-educated person to befriend an educated person. Green et al. (1985) and Basu and Foster (1998) argue that the benefits of education are shared *within* the same household.

For example, a literate family member may help a non-literate member to fill in a form or read a brochure left behind by the extension officer. Along the same line of thinking, households who have no literate members may find it interesting to befriend households with literate members. We create two dummy variables to test for the effect of education on link formation. The first dummy takes on a value of one when one household has no members who have completed primary and the other household has at least one member who has completed primary education. The second dummy indicates that neither of the two households have any household members who have completed primary education. The reference is the category in which both households have at least one educated member. Of course we have to bear in mind that education can proxy for other variables, like wealth.

2.7 Activity overlap

Households who are engaged in similar activities are likely to have covariate income streams, which we expect to be an impediment to link formation. To test for this we create an index that measures how similar the income portfolio of two households is. All the income generating activities the household engages in, were placed into 7 categories (casual labour, trade, crop production, livestock rearing, assets and processing farm produce, other off-farm activities). Each household can have several activities within the same category. By dividing the number of activities in each category by the total number of activities the household is engaged in, we get a rough measure of the spread of activities expressed in percentages of the total portfolio (it is only rough because we take no account of the income each activity generates). The activity overlap between two households is taken to be the sum of the minimum percentages in each category. E.g. if the portfolio of A is (30%, 0%, 40%, 0%, 30%, 0%, 0%) and of B is (0%, 20%, 70%, 0%, 10%, 0%, 0%) then the overlap will be $0\% + 0\% + 40\% + 0\% + 10\% + 0\% + 0\% = 50\%$. The overlap will always lie between 0% and 100%. In our sample the mean and median activity overlap are both 0.55 and the standard deviation is 0.25.

2.8 Difference in the age of the household heads

We want to know whether insurance partners are chosen within one's own age category, outside one's own age category, or whether age doesn't matter. The neediness of a household and its possibilities to reciprocate in the future might depend on its age. We use the age of the household head to proxy for the age of the household. A possible measure of how heterogeneous two households are with respect to age is the age difference between the two heads.

2.9 Differences in wealth

We would like to know whether networks form between households of similar wealth. If so, we expect that the wealth differences between two households have a negative effect on link formation. Table 5 shows the distribution of livestock and land differences across the dyads.

3 Econometric issues

This section first describes the three different sets of regressions that were performed on the data and then goes on to discuss possible problems of autocorrelation. In a first set of regression we only enter relational variables. Still, attribute variables might play an important role in explaining the degree of connection between two households. For example, older households have had more time to build up links, so we might expect them to have more links than younger households. In this case not only the difference in age between the two households matters, but also the absolute age of both parties.

To control for all attribute variable we run a second set of regression in which we control for household fixed effects. To this end we introduce a special kind of dummy variable. There are 119 of these dummies in total, one for each household in the sample, and each of them indicates whether that specific household is part of the pair or not. This means that every row of the data always contains two dummy variables equal to one (not taking account of the reference category). On the one hand, we control for observable attribute variables. For example, we will find that richer households typically have more links than poorer households. On the other hand, by including dummies we will also control for unobserved attribute variables. For example, a very communicative and cheerful person might be likely to have more links than an introvert. Related, how many network partners you choose, might be correlated with unobservables. Once dummies are included in the regression we have controlled for these effects.

We also run a third set of regressions in which, instead of entering dummies, we enter the maximum age of the household head, livestock value, land ownership and the number of respondents. This has essentially the same effect as the inclusion of the dummies. Now we do not perfectly control for all non-relational variables, but we do get information on which non-relational effects are at work. The number of respondents is entered to control for the fact that aggregating individual interviews to household level causes households with many respondents to have more links. The mean and median number of respondents is 2. The distribution of the maximum values across dyads is summarized in Table 5.

Before reporting the regression results, we want to draw attention to the issue of autocorrelation. Say u_{ij} is the error term associated with the dyad formed by household i and household j . We should be particularly wary of possible correlation between u_{ij} and all u_i , u_j and u_k (i.e. all other dyads in which i or j appear). The unobserved attributes of i feature in the error terms of all the dyads containing i . For example, if household i consists of grumpy, ill-tempered individuals, then all other households may avoid having insurance links with them for reasons beyond anything we observe in the data. This would make all u_{ij} and u_i correlated. Thus in the first set of regressions, which has no controls for attributes, autocorrelation may bias our results. In the second set of regressions, household fixed effects purges out the effects of all attribute variables and therefore eliminates the autocorrelation in the example.⁴ In the third set of regressions only observed attribute variables are controlled for. To be sure, even in the household fixed effects regressions we need to make an assumption about the error

⁴ In the literature on network analysis this autocorrelation problem has long been recognized and is solved by running QAP regression in stead of using dummies (e.g. Krackhardt 1988).

structure: the error terms of two dyads, which contain no mutual members are assumed to be uncorrelated to each other.

4 Regression results

Tables 6–8 report the three different sets of regression results. The coefficients of the logit regressions in the third column of each table are marginal effects, so we can interpret them as the increase in probability of a link after a unitary increase in the explanatory variable at the sample mean. Throughout all the different regression specifications there are four variables which stand out because of the consistency of their sign and their high significance. These are the kinship dummies, the geographical distance, being of the same religious affiliation and the number of common friends.

Kinship has the strongest effect on link formation. Compared to having no blood bond, sharing 50 per cent of one's genetic material (parents, children and siblings) raises the chances of having a link with 39 per cent in the household fixed effects logit regressions. As the genetic distance increases, the effect diminishes.

Increasing the distance between two households with 1 km, reduces the probability of having a link with approximately 9 per cent. Enforcement constraints seem to outweigh issues related to the non-covariance of income here. To put these numbers into perspective, remember that the two furthest neighbours in the village live 1.7 km away from each other.

It is common practice in economics to use the term 'the network of family, neighbours and friends'. Our results confirm that family and neighbours do indeed go far in explaining group formation. 'Friends', I believe, is supposed to be the rest-term for all who are neither neighbour nor kin, but still network members. The remainder of this section studies what the network looks like once kinship and geographical proximity are controlled for.

Being of the same religious affiliation has a significant, but small effect on network formation. Clan membership has no influence in the regressions which concentrate solely on direct links. This may be because there are 25 clans amongst 119 households, which leaves few possibilities for direct matching along clan lines. It is perhaps not surprising that the variable does become significant in the geodesic distance regression. Once we make more comprehensive use of the network graph, by also considering what lies *beyond the direct links* clan does become significant; i.e. clan mates do lie closer to each other on the network graph. However, we should be cautious about this result as the significance disappears in the fixed effects regressions.

Having many common friends increases the possibility that two households are also linked. The size of this effect is smaller in the household fixed effects regression than in the others, which may indicate that it is correlated to unobserved fixed effects.

Both education dummies have negative coefficients. Although the effects only become significant in the geodesic distance regressions, the sign of the first education dummy does seem to suggest that households without any educated members avoid having insurance links with each other. At the same, the second education dummy tells us that dyads that are mixed in terms of education also have less chance of being linked to each

other. This implies that inter-household externalities to education do not play a role in network formation. On the contrary, households with educated members seem to lie closer to each other on the network graph. We will find a similar result with respect to wealth.

The variable ‘activity overlap’ is significantly positive in the geodesic distance regressions, which suggests that non-covariance of income is less of an issue than information constraints here. The effect becomes smaller in the fixed effects regressions. The most likely cause of this difference is correlation between the fixed effects and the activity overlap. Households with many activities are likely to have a high activity overlap. At the same time, the regressions in Appendix 2 show that households with many activities have many links. Thus in the regressions without attribute controls the activity overlap picks up some of the effect of the number of activities. Indeed, running the third set of regressions (attribute controls through maximum values) with a new variable ‘maximum number of activities’ included, gives results which are very similar to the fixed effects results (results are not shown). Still, it is worrying that in the regressions for direct links the effect of activity overlap is no longer significant. One reason might be that we have defined the income categories too broadly and that even within one category there might be very non-correlated income streams.

Age differences between households have a very small, significant effect. There is a slight tendency for households to choose network partners close to their own age. In the third set of regressions, we see that the maximum age of the household has a small, but significant effect on link formation. This might be because older households have, through the years, established more links than younger households have.

The differences across households in livestock value and land holdings, are measures of how heterogeneous a pair of households is with respect to wealth. The coefficients are positive in the regressions without attribute controls. This would suggest that larger disparities in wealth enhance network formation and the network is redistributive. This is, however, not what is really going on. Fafchamps (1992) notes that wealthier people are more desirable to befriend and this is also apparent in the Nyakatoke data. The total number of links a household has correlates with its livestock and land holdings (regression results are given in Appendix 2, Table A2.1): the rich have denser networks than the poor. This popularity effect implies that a dyad with a rich person in it has more chance of being linked. At the same time, such a dyad will also have a large wealth disparity. This is obviously so if the other half of the dyad is a poor person, but is also true in combination with another rich person, because of scale effects –in absolute terms the rich have bigger wealth gaps between each other than the poor.

Including household dummies purges the regression of household fixed effects, so also any popularity effect. We can see that all the wealth variables get reverse signs once dummies are included. Including the maximum (across the two households) livestock value or landholdings gives the same results, as they will also pick up any popularity effect of wealthy households. The picture we then get is that rich households choose each other as network partners, but poor households avoid each other as network partners.

Table 6
Regressions excluding attribute controls

	Total No. of links reported by individuals (ordered logit)		Total No. of household level links (ordered logit)		Dummy (logit*)		Geodesic distance** (ordered logit)	
	coeff.	p	coeff.	p	coeff.	p	coeff.	p
Parents, children and siblings	2.502	0.00	2.371	0.00	0.346	0.00	2.921	0.00
Nephews, nieces, uncles, aunts, cousins, grand-parents and grandchildren	1.497	0.00	1.468	0.00	0.129	0.00	1.132	0.00
Other blood bond	1.182	0.00	1.230	0.00	0.083	0.00	0.486	0.01
Distance (km)	-1.999	0.00	-2.003	0.00	-0.090	0.00	-0.763	0.00
No. of common friends	0.445	0.00	0.444	0.00	0.020	0.00	2.059	0.00
Same religion	0.344	0.00	0.336	0.00	0.016	0.00	0.147	0.00
Same clan	0.146	0.37	0.120	0.46	0.006	0.41	0.215	0.02
Neither of the 2 HHs has a member who completed primary	-0.337	0.24	-0.413	0.15	-0.015	0.17	-0.717	0.00
Only 1 HH has a member who completed primary	-0.167	0.14	-0.167	0.14	-0.007	0.14	-0.313	0.00
Activity overlap	-0.065	0.80	-0.063	0.81	-0.003	0.77	0.695	0.00
Difference in age of HH head (10 years)	-0.087	0.03	-0.098	0.01	-0.004	0.02	0.005	0.79
Difference in livestock value (/100000 TSh)	0.044	0.01	0.040	0.02	0.002	0.03	0.043	0.00
Difference in land (Ha)	0.061	0.12	0.050	0.20	0.002	0.16	0.001	0.95
pseudo R ²	0.16		0.17		0.20		0.24	
p-value of chi ² test	0.00		0.00		0.00		0.00	
N	6555		6555		6555		6555	

* Marginal effects at the sample means.

** The coefficients in the ordered logit for geodesic distance have been multiplied by -1. This makes it easier to compare them to the coefficients from other regressions. Thus, a positive coefficient here means that the variable REDUCES the geodesic distance, i.e. they are better connected.

Source: Nyakatoke Household Survey.

Table 7
Regressions including attribute controls: household fixed effects

	Total No. of links reported by individuals (ordered logit)		Total No. of household level links (ordered logit)		Dummy (logit*)		Geodesic distance** (ordered logit)	
	coeff.	p	coeff.	p	coeff.	p	coeff.	p
Parents, children and siblings	2.880	0.00	2.739	0.00	0.391	0.00	3.277	0.00
Nephews, nieces, uncles, aunts, cousins, grand-parents and grandchildren	1.749	0.00	1.694	0.00	0.143	0.00	1.053	0.00
Other blood bond	1.387	0.00	1.426	0.00	0.087	0.00	0.560	0.00
Distance (km)	-2.641	0.00	-2.627	0.00	-0.092	0.00	-1.174	0.00
No. of common friends	0.210	0.00	0.203	0.00	0.008	0.00	1.775	0.00
Same religion	0.408	0.00	0.401	0.00	0.015	0.00	0.183	0.00
Same clan	0.130	0.46	0.081	0.65	0.003	0.70	0.026	0.80
Neither of the 2 HHs has a member who completed primary	-1.216	0.32	-1.239	0.31	-0.026	0.31	-3.652	0.00
Only 1 HH has a member who completed primary	-0.613	0.31	-0.590	0.33	-0.019	0.33	-1.739	0.00
Activity overlap	-0.344	0.36	-0.321	0.39	-0.013	0.33	0.291	0.08
Difference in age of HH head (10 years)	-0.151	0.00	-0.164	0.00	-0.005	0.00	-0.027	0.27
Difference in livestock value (/100000 TSh)	-0.132	0.06	-0.169	0.02	-0.008	0.01	-0.051	0.51
Difference in land (Ha)	-0.088	0.21	-0.095	0.18	-0.003	0.16	-0.025	0.52
Pseudo R ²	0.20		0.22		0.25		0.34	
p-value of chi ² test	0.00		0.00		0.00		0.00	
N	6555		6555		6555		6555	

* Marginal effects at the sample means.

** The coefficients in the ordered logit for geodesic distance have been multiplied by -1. This makes it easier to compare them to the coefficients from other regressions. Thus, a positive coefficient here means that the variable REDUCES the geodesic distance, i.e. they are better connected.

Source: Nyakatoke Household Survey.

Table 8
Regressions including attribute controls: maximum values

	Total No. of links reported by individuals (ordered logit)		Total No. of household level links (ordered logit)		Dummy (logit*)		Geodesic distance** (ordered logit)	
	coeff.	p	coeff.	p	coeff.	p	coeff.	p
Parents, children and siblings	2.573	0.00	2.444	0.00	0.352	0.00	2.955	0.00
Nephews, nieces, uncles, aunts, cousins, grand-parents and grandchildren	1.549	0.00	1.517	0.00	0.135	0.00	1.150	0.00
Other blood bond	1.108	0.00	1.165	0.00	0.074	0.00	0.439	0.01
Distance (km)	-2.111	0.00	-2.114	0.00	-0.091	0.00	-0.790	0.00
No. of common friends	0.392	0.00	0.391	0.00	0.018	0.00	2.035	0.00
Same religion	0.357	0.00	0.350	0.00	0.016	0.00	0.153	0.00
Same clan	0.167	0.31	0.134	0.42	0.007	0.37	0.232	0.01
Neither of the 2 HHs has a member who completed primary	-0.194	0.52	-0.288	0.34	-0.009	0.42	-0.588	0.00
Only 1 HH has a member who completed primary	-0.102	0.41	-0.113	0.36	-0.004	0.46	-0.221	0.00
Activity overlap	-0.087	0.74	-0.077	0.76	-0.004	0.70	0.689	0.00
Difference in age of HH head (10 years)	-0.173	0.00	-0.194	0.00	-0.008	0.00	-0.054	0.05
Difference in livestock value (/100000 TSh)	-0.182	0.14	-0.232	0.08	-0.015	0.02	-0.174	0.20
Difference in land (Ha)	-0.286	0.00	-0.284	0.00	-0.012	0.00	-0.280	0.00
Max. age of HH head (10 years)	0.124	0.01	0.134	0.01	0.005	0.02	0.078	0.00
Max. livestock value (/100000 TSh)	0.205	0.09	0.252	0.05	0.016	0.01	0.201	0.14
Max. land (Ha.)	0.320	0.00	0.309	0.00	0.013	0.00	0.257	0.00
Max. No. of respondents	0.132	0.13	0.109	0.21	0.006	0.13	0.137	0.00
Pseudo R ²	0.17		0.18		0.21		0.25	
p-value of chi ² test	0.00		0.00		0.00		0.00	
N	6555		6555		6555		6555	

* Marginal effects at the sample means.

** The coefficients in the ordered logit for geodesic distance have been multiplied by -1. This makes it easier to compare them to the coefficients from other regressions. Thus, a positive coefficient here means that the variable REDUCES the geodesic distance, i.e. they are better connected.

Source: Nyakatoke Household Survey

5 Concluding remarks

On the basis of household and network data collected in a Haya village in rural Tanzania, we found that kinship, geographical proximity, the number of common friends, clan membership, religious affiliation and wealth strongly determine network formation.

We would like to point to some of the shortcomings of our analysis. First, we have excluded links to households living outside Nyakatoke and it is difficult to determine how this influences the results. Second, we conduct the analysis at household level, so we abstract from any intra-household issues. Third, the identification of which we call a network partner may depend on the framing of the question.

Insights in endogenous network formation are important for assessing vulnerability of households. A vulnerability assessment should distinguish between households that are likely to experience network shocks (everyone in their network is hit at the same time) and those who are not. Note that a shock might be common in the sense that everyone in the economy is hit (e.g. harvest failure in a large area), but still it can affect different networks in a different way. We might find that there are weak networks that collapse under this shock and strong networks that can cope with the shock. This is likely if households tend to link up with others of similar wealth, occupation and place of residence.

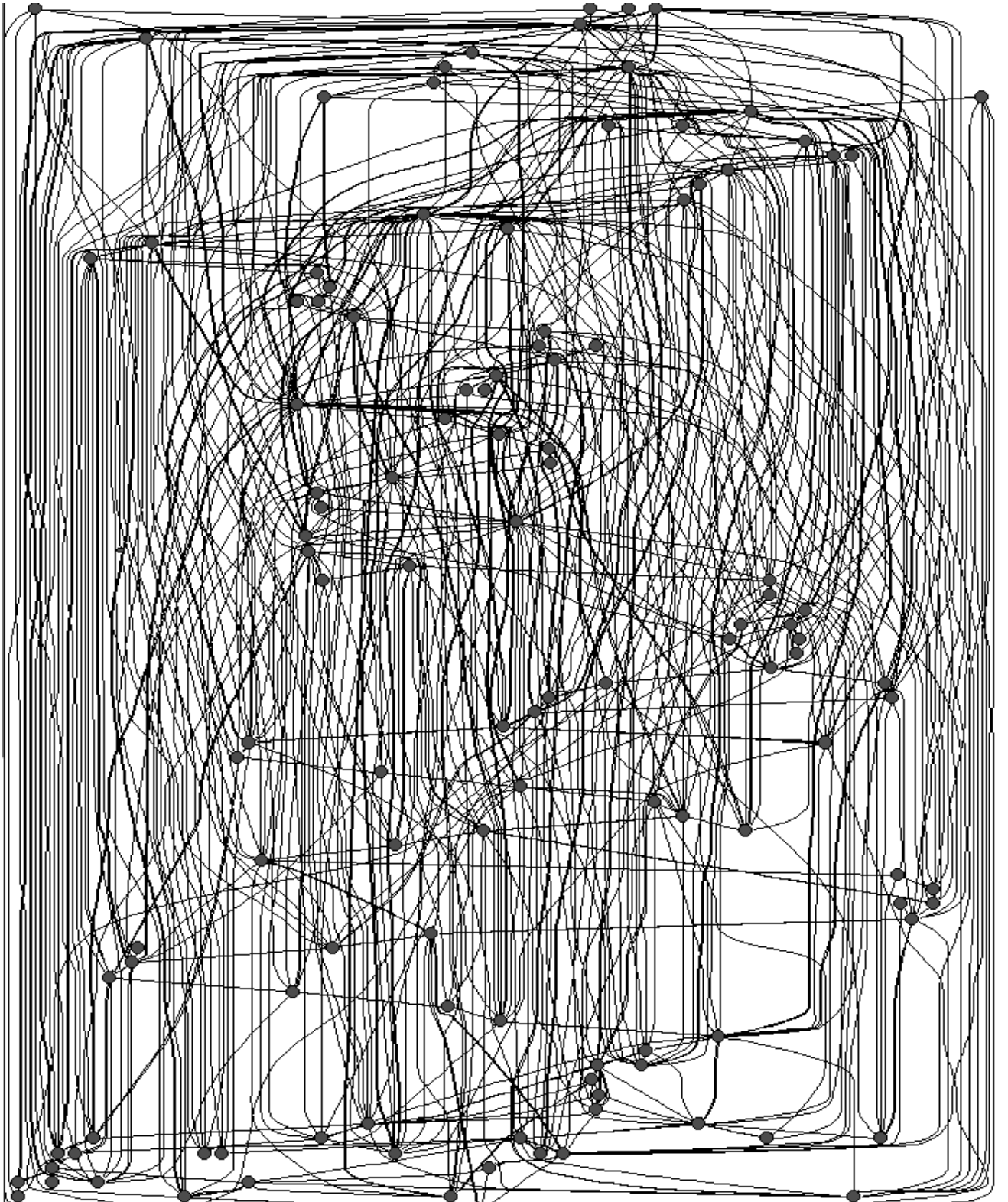
Better insight in the determinants of group formation might also point to categories of vulnerable households who fail to enter into risk-sharing arrangements. Our data suggest that poor households have less dense networks than the rich, making them more vulnerable in the face of idiosyncratic risk.

References

- Bala, V. and Goyal, S. (2000), 'A Non-cooperative Model of Network Formation', *Econometrica*, Vol.68, No.5, pp.1181–229.
- Basu, K. and Foster, J. (1998), 'On Measuring Literacy', *Economic Journal*, Vol.108 (451), pp.1733–49.
- Coate, S. and Ravallion, M. (1993), 'Reciprocity without Commitment: Characterisations and Performance of Informal Risk-sharing Arrangements', *Journal of Development Economics*, Vol.40, pp.1–24.
- Fafchamps, M. (1992), 'Solidarity Networks in Preindustrial Societies: Rational Peasants with a Moral Economy', *Economic Development and Cultural Change*, Vol.41, No.1, pp.147–74.
- Fafchamps, M. and Lund, S. (2000), '*Risk-sharing Networks in Rural Philippines*', mimeo, Oxford: Oxford University.
- Genicot, G. and Ray, D. (2000), 'Endogenous Group Formation in Risk-Sharing Arrangements', mimeo, University of California and New York University.
- Goldstein, M. (2000), 'Intra-household Allocation and Farming in Southern Ghana', Ph.D. dissertation, University of California at Berkeley.
- Green, S., Rich, T. and Nesram, E. (1985), 'Beyond Individual Literacy: The Role of Shared Literacy for innovation in Guatemala', *Human Organization*, Vol.44, pp.313–21.
- Grimard F. (1997), 'Household Consumption Smoothing through Ethnic Ties: Evidence from Côte d'Ivoire', *Journal of Development Economics*, Vol.53, pp.391–421.
- Krackhardt, D. (1988), 'Predicting with Networks: Nonparametric Multiple Regression Analysis of Dyadic Data', *Social Networks*, Vol.10, pp.359–81.
- Mitti, R. and Rweyemamu, T. (2000), 'Taswira ya Kijamii na Kiuchumi ya Kitongoji Nyakatoke', mimeo, KU Leuven.
- Murgai, R., Winters, P., Sadoulet, E. and de Janvry, A. (2000), 'Localized and Incomplete Mutual Insurance', mimeo, World Bank and University of New England and University of California.
- Platteau, J.P. (1991), 'Traditional Systems of Social Security and Hunger Insurance: Past Achievements and Modern Challenges' in Ahmad, E., Drèze, J., Hills, J., and Sen, A. (eds), *Social Security in Developing Countries*, Oxford: Clarendon Press, pp.112–70.
- Reining, P. (1967), 'The Haya: the Agrarian System of a Sedentary People', Ph.D. dissertation, University of Chicago.
- Spinnewyn, F. and Wijaya, M. (1996), 'Voluntary Reciprocity and Income Mobility', CES discussion paper, KU Leuven.
- Wasserman, S. and Faust, K. (1995), '*Social Network Analysis*', Cambridge: Cambridge University Press.

Appendix 1

The network graph of Nyakatoke. Each dot is a household, each line a link. Households are *not* positioned according to geographical location. The map has been drawn using a programme called DOTTY from Graphviz.



Appendix 2

OLS for the total number of household links.

Table A2.1
OLS for the total number of household links (N=116)

	coeff.	p-value of t-stat	coeff.	p-value of t-stat	coeff.	p-value of t-stat
Constant	-1.037	0.56	0.443	0.80	0.502	0.77
No. of activities	0.376	0.02	0.554	0.00	0.539	0.00
One member completed primary	0.779	0.44	0.922	0.38	1.507	0.13
Age HH head	0.059	0.02	0.063	0.02	0.069	0.01
Livestock*	0.652	0.00			0.760	0.00
Land	0.134	0.72	0.811	0.03		
No. of respondents	1.771	0.01				
R squared	0.35		0.23		0.28	
p-value of F-stat.	0.00		0.00		0.00	

* coefficient multiplied by 100000