



UNIVERSITY OF
CAMBRIDGE

Cambridge Working Papers in Economics

Does Anti-Diversification Pay? A One-Sided Matching Model
of Microcredit

Thilo Klein

June 19, 2015

CWPE 1521

Does Anti-Diversification Pay? A One-Sided Matching Model of Microcredit*

Thilo Klein[†]

July 19, 2015

Abstract

In many economic situations, market participation requires that agents form groups subject to exogenous rules. Consider a microfinance institution that decides on rules for diversifying borrower groups in terms of their exposure to income shocks. Such rules affect group repayment by influencing both who matches with whom (direct effect) and who participates in the market (participation). I develop the key trade-off for conflicting predictions of extant theoretical models and estimate both effects separately. Group formation creates an endogeneity problem, but a matching model exploits the exogenous variation from counterfactual groups. I find that while diversification has no participation effect it has a significant positive direct effect.

Keywords: microcredit; joint liability; risk; diversification; market design; stable matching; endogeneity; selection model; agriculture; Thailand

JEL Classifications: C11, C31, C34, C36, C78, C57, D02, D47, D82, G21, O16, Q14

1 Introduction

Economists often observe data on interactions: who interacts with whom, which students team up in study groups, which characters form entrepreneurial teams and which firms merge with each other. The empirical analysis of the outcomes

*I thank my advisor Paul Kattuman and my examiners Maitreesh Ghatak and Hamish Low for their guidance. I thank Christian Ahlin, Kumar Aniket, Britta Augsburg, Peter Burgold, Michael Freeman, Susann Giese, Michael Goller, Pramila Krishnan, Robert Lensink, Vincent Mak, Stefan Scholtes, Susan Steiner, seminar participants at Cambridge, Kristiansand and Geneva for comments. The research reported herein was supported by the Cambridge Home and EU Scholarship Scheme (CHESS), German Academic Exchange Service (DAAD) and Economic and Social Research Council (ESRC). All views and errors are mine.

[†]Organisation for Economic Co-operation and Development, 2 rue Andre Pascal, 75775 Paris Cedex 16, France. Email: thilo.klein@oecd.org

of such interactions is complicated by three, possibly counteracting, effects, which this paper distinguishes empirically using credible exclusion restrictions. To illustrate, take the decision of a firm's top management on the optimal intercultural mix of the firm's units.

- First, there is the *direct effect* of intercultural team composition on team outcomes for a given set of workers. This concerns, for instance, whether communication problems outweigh the synergies within mixed teams.
- Second, this direct effect is net of *sorting bias*. This bias arises if, for example, open-minded workers are more likely to sort into mixed teams and open-mindedness (i) is unobserved by the econometrician and (ii) results in better team outcomes. In this case, the direct effect of mixed teams would be overstated because it picks up the positive effect of open-mindedness, which is unobserved.
- Third, a management decision stipulating mixed teams would have a *participation effect*, in that it may result in a smaller pool of applicants and workers resigning if they dislike working in mixed teams.

The development economics literature on group lending in microfinance provides a compelling case for the study of these effects. Microfinance has pushed out financial frontiers in developing countries in terms of expanding access to credit for low-income households that lack seizable collateral. This has been enabled by a critical innovation in contract design, namely joint liability, or group lending. This contract form is both the most relevant in the field and the most studied in economic research on microcredit. In group lending, borrowers form groups endogenously. This, the economic literature demonstrates, is socially preferable when matching is based solely on the borrowers' risk type. The implications of matching based on other dimensions are less clear, however. When matching also has adverse effects on repayment, banks would be well advised to impose restrictions on permissible group constellations. In practice, banks operating on the Grameen model explicitly rule out the grouping of relatives in order to avoid collusion against the lender (see [Alam and Getubig, 2010](#), p. 17). An understanding of how group composition affects repayment is therefore of very practical importance for banks as well as for our understanding of the economic theory of joint-liability lending.

The focus of this paper is on the effect of matching on exposure to similar income shocks, as is common in agriculture. In this context, both academics and

microfinance institutions – most famously the Grameen Bank – have questioned the use of group lending because group members avoid joint-liability payments when their projects fail simultaneously (see [Ghatak, 2000](#)). Contrary to this widely held view, [Ahlin and Townsend \(2007\)](#) extend two well-established models of joint-liability lending to show that positive project covariation can raise repayment. The first, the [Ghatak \(1999\)](#) model, considers an adverse selection setting á la [Stiglitz and Weiss \(1981\)](#) where a pooling contract subsidises risky projects. Groups with socially productive, safe risk types therefore find it less beneficial to borrow and are drawn out of the market. The second is the ([Stiglitz, 1990](#)) ex-ante moral hazard model, in which group members choose cooperatively between safe and risky projects. Safe projects are always preferable socially but not necessarily privately. In both models, project correlation makes borrowers avoid liability payments. This is because it makes it less likely that one borrower’s project will succeed while her partner defaults. This, the authors show, has a positive *participation effect*, in that it draws safe types back into the market in [Ghatak \(1999\)](#), and also a positive *direct effect*, in that it makes choosing safer projects more attractive in [Stiglitz \(1990\)](#). They confirm these predictions using data from Thailand. Repayment implications and empirical results are summarised in Table 1, Columns 1-i and 1-ii.

Table 1: Summary of theoretical and empirical results

Upward arrows indicate a positive repayment effect of project covariation.						
	Ahlin/Townsend		Theory	This paper		
	Theory	Empirics		Empirics		
		Logit		Probit	Structural	Simulation
	(1-i)	(1-ii)	(2-i)	(2-ii)	(2-iii)	(2-iv)
A. Direct effect						
- Stiglitz (1990)	↑		↓ ^{a)}		↓	↓ ^{c)}
B. Participation						
- Ghatak (1999)	↑		↓ ^{b)}			↓
<i>subtotal</i> (A+B)	↑		↓			↓
C. Sorting bias			↑		↑	
<i>total</i> (A+B+C)		↑		↑	↑	

^{a)} Negative for low risk aversion *or* low return differential of risky and safe projects.

^{b)} Negative for low marginal risk types *or* high liability payment.

^{c)} Coefficients for simulations are taken from the estimates of the structural model.

The contributions of this paper are threefold. First, it develops the key trade-off of the conflicting effects suggested separately in the literature. The negative effect of anti-diversification in Ghatak (2000) is found to be dominant for a wide range of parameter constellations (see Table 1, 2-i).

Second, it establishes uniqueness conditions for equilibria in group formation games with non-transferable utility in finite markets. A unique equilibrium is necessary for the likelihood of empirical models to be well-defined. In the theoretical framework, matching is on two dimensions: borrowers' risk type and which of two external shocks borrowers' projects are exposed to. Preferences are *aligned* in risk type – agents prefer safer partners who are less likely to need bailing out – and *assortative* in exposure type – agents prefer group members of the same exposure type because they are more likely to fail simultaneously and thereby avoid joint-liability payments. Such preference profiles are also common in other contexts. In marriage markets (Banerjee *et al.*, 2013), partners have been shown to marry up within the same caste. Similarly, in study groups, pupils match with more academically able peers within the same gender. In agricultural lending, there is one dominant exposure type: exposure to weather shocks. If these shocks are strong, then matching is first on exposure type and then aligned in risk type (within exposure type). In villages with two lending groups, this results in one *dominant group*, which is composed of the 'weather shock' exposure types with the safest projects. This group is in equilibrium because no borrower would prefer to match with a different exposure type or a worse risk type. The remaining borrowers form a *residual group*, which is composed of 'weather shock' exposure types with riskier projects and those with other exposure types.

Third, the paper develops a structural model that corrects for *sorting bias*; implemented in R package `matchingMarkets` (Klein, 2015b). This bias arises because in equilibrium the dominant group has, on average, projects that are both safer and more highly correlated than the residual group. If risk type is partially unobservable and captured in the error term, then the repayment effect of project correlation is biased upwards (Table 1, 2-i). The structural model is similar to the Heckman (1979) selection correction but generalises this model to allow for the selection process being the equilibrium outcome of a group formation game. The identifying exclusion restriction for the direct effect is that the characteristics of all agents in the market affect who matches with whom, but the performance of a matched group is determined only by its own members. No additional instrumental variables are required. This is crucial in a context where instruments are impossible to find, because sorting occurs with a view to optimise group outcomes.

The model is applied to a resurvey of the data used in [Ahlin and Townsend \(2007\)](#), allowing me to model the endogenous group formation explicitly. The *direct effect* is a test of the revised predictions in the moral hazard model of [Stiglitz \(1990\)](#). In line with the predictions, I find a significantly negative direct effect of project covariation on repayment. This negative direct effect is net of a positive *sorting bias* that – if not controlled for – would yield an erroneous, positive estimate of the direct effect (Table 1, 2-iii). The *participation effect* from matching on risk exposure is tested in agent-based simulations using parameter estimates from the structural model. Varying the rules to either allow or prohibit matching on risk exposure – but keeping the model fixed to predict the repayment outcome – allows me to separate the participation effect in my revision of [Ghatak \(1999\)](#). I find that anti-diversification draws borrowers into the programme who would not have taken a loan otherwise. However, as predicted by the model, this positive effect is more than offset by the negative effect of the bank losing joint-liability payments when projects fail simultaneously (see Table 1, 2-iv).

Taken together, this paper reconciles predictions from the theoretical models and empirical evidence from Thailand with the literature, which has long considered the positive project covariation of agricultural loans as an impediment to the expansion of lending programmes in rural markets (refer to [Mosley, 1986](#)). The results also add empirical evidence to the long-running discussion on group versus individual lending. They suggest that joint liability is less effective in agricultural lending and that this adverse effect is exacerbated by endogenous matching on common risk exposure. In contexts where joint-liability contracts are desirable, lenders should prevent the grouping of borrowers who are exposed to similar income shocks.

This paper is organised as follows. Section 2 reviews the literature. Section 3 develops the key trade-off between the conflicting effects suggested in the literature and establishes uniqueness conditions for equilibrium matching in non-finite markets. Section 4 presents the empirical strategy. Section 5 describes the data and presents the results. Section 6 concludes.

2 Literature review

There are two parts to this literature review. The first provides a map of the issues covered in the paper and signposts to connect them to broader topics. The second presents theoretical models of repayment in joint-liability lending and the revisions I make to them in the theoretical framework in Section 3.

2.1 Empirical work

Extant empirical work on correlated returns in joint-liability lending has produced ambiguous results: of four key studies, three find a significantly negative repayment effect and one – the most recent – finds the opposite. [Wydick \(1999\)](#), [Zeller \(1998\)](#) and [Sharma and Zeller \(1997\)](#) follow a common methodology in that they measure the covariation of project returns using occupational heterogeneity within groups. This measure is problematic because it also captures the cost of monitoring between group members. [Ahlin and Townsend \(2007\)](#) construct a measure of common exogenous shocks within groups: the probability that the worst business year in a five-year time window coincides for two randomly chosen group members. This is the measure I use.

Lab and field experiments

The findings of previous studies are subject to a sorting bias that is well recognised in the literature (see [Hermes and Lensink, 2007](#)). To overcome this issue, experimental methods have become popular means of testing theories of joint-liability lending. [Karlan \(2007\)](#) makes use of the quasi-random group assignment of microlender FINCA in Peru to estimate the *direct effect* of social connections. In framed field experiments, [Giné et al. \(2010\)](#) implement a ‘partner choice’ treatment to estimate the *direct effect* of endogenous group formation compared to random assignment. Similarly, *participation* and *direct effects* combined can be tested with ‘group recruitment’ ([Abbink et al., 2006](#)) or ‘self-selection’ ([Cassar and Wydick, 2010](#)) treatments that require participants to register for lab experiments in groups. The main advantage of the technique I develop in this paper is that it can be applied to field data to test the effects of sorting on specific, policy-relevant variables rather than the effect of sorting per se. This technique thus allows me to derive concrete recommendations for designing rules for group formation.

Networks in microfinance

This paper is also connected with the broader issues of networks and network formation in microfinance, which have gained much attention recently. In the context of informal village networks, [Fafchamps and Gubert \(2007\)](#) use dyadic regression to identify the formation and determinants of risk sharing in informal insurance networks among villagers. The main difference between these informal and formal networks (such as the groups studied in this paper) is that the former have no restrictions on group size. Restricting group size results in competition for

places, which both creates the interactions that complicate the empirical analysis of these markets and provides a valuable source of exogenous variation. From this viewpoint, the work that is most closely related to mine is that of [Klonner \(2006\)](#) and [Eeckhout and Munshi \(2010\)](#) on the matching of fixed-size chit fund groups. The main difference, however, is that the market for chit funds is two-sided in that it brings together borrowers and lenders. By contrast, I study one-sided matching, where anyone can match with anybody else.

Structural empirical work using matching models

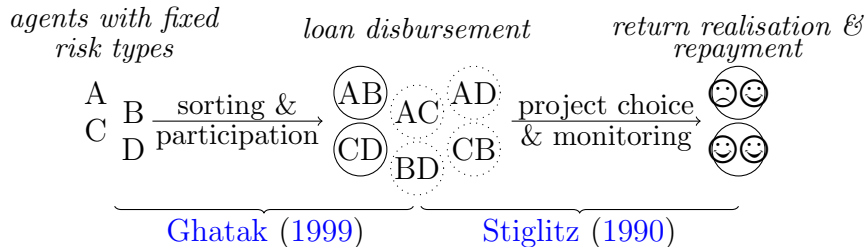
Structural empirical work on matching markets has a wide range of beneficial applications. In one-sided matching, these applications range from US school district mergers ([Gordon and Knight, 2009](#)) to Japanese municipal amalgamations ([Weese, 2015](#)), which are modelled as stable roommates and group formation games, respectively, in order to understand the determinants of mergers. The work most closely related to my paper is the analysis of microfinance group formation in [Ahlin \(2009\)](#). The value added by my paper is that I also analyse the implications of endogenous group formation for repayment. I thereby add to pioneering work on selection models in matching markets, which simultaneously estimate a matching model that parametrically selection-corrects an outcome equation. Such models have been proposed for two-sided markets in [Sørensen \(2007\)](#), who examines whether firms are more likely go public when matched with more experienced venture capitalists, and applied in [Chen \(2013\)](#) and [Park \(2013\)](#). The model developed here is the first to implement this strategy in a one-sided matching market.

2.2 Theoretical models of joint-liability lending

In moving from empirics to theory, it should be noted that this paper is not concerned with finding the optimal credit contract. The focus is instead on optimal market design (see [Roth, 2008](#), for an overview), i.e. how to set the rules for group formation such that group repayment is maximised, taking the contract terms as given. I use terms from a joint-liability contract because this is demonstrably the optimal contract form in a setting with correlated returns and a lack of seizable collateral. In a moral hazard (effort choice) context, [Che and Yoo \(2001\)](#) find that joint liability is the optimal collusion-proof contract, even under almost perfect project correlation. Similarly, in the [Stiglitz and Weiss \(1981\)](#) adverse selection setting, [Laffont \(2003\)](#) finds that joint-liability lending is still the optimal contract when returns are correlated.

The theoretical framework builds on the adverse selection model of [Ghatak \(1999\)](#) and the ex-ante moral hazard model of [Stiglitz \(1990\)](#).¹ Figure 1 illustrates the sequence of events in group lending. [Ghatak \(1999\)](#) models who matches with whom (sorting) and who participates in the market, given their group members. A key result of this model is the homogeneous matching of agents into equilibrium groups based on risk type. For this model, [Ahlin and Townsend \(2007\)](#) find that correlated returns improve repayment because they make borrowing more attractive and thereby draw safe types (who would not have borrowed otherwise) into the market. I show that this positive effect is dominated by a negative anti-diversification effect from banks losing joint-liability payments when projects fail simultaneously. After loan disbursement, agents decide together whether to gamble and realise riskier projects ([Stiglitz, 1990](#)). This decision depends on both agent and project characteristics. [Ahlin and Townsend \(2007\)](#) argue that positive project correlation lowers the temptation to gamble and improves repayment. I show that this is only the case if (i) borrowers are strongly risk averse *and* (ii) the risky project has considerably higher returns than the safe project.

Figure 1: Group lending (sequence of events).



3 Theoretical framework

This section is divided into two subsections. The first derives the repayment implications of correlated project returns for the two most widely cited theoretical models of joint-liability lending. The second (i) establishes uniqueness conditions for the equilibrium matching used in the empirical model and (ii) derives the sign of the sorting bias that results when matching is on both risk and exposure type.

¹While attention is restricted to these two models, other effects, such as ex-post moral hazard, may be at work.

3.1 Revised theories and implications

For both theoretical models I present the model and the positive repayment effects derived in the model extensions by [Ahlin and Townsend \(2007\)](#). I then introduce the negative effect of anti-diversification in [Ghatak \(2000\)](#) and develop the key trade-off.

3.1.1 Adverse selection: Participation effect

The [Ghatak \(1999\)](#) model uses the [Stiglitz and Weiss \(1981\)](#) setting of credit rationing. There is a continuum of risk-neutral borrowers who are endowed with one unit of labour and no pledgeable collateral. Agents can either sell their labour and earn an outside option \bar{u} or borrow and invest one monetary unit in an uncertain project. Agent i 's project yields an actual outcome of y_i with success probability p_i and 0 otherwise. The distribution of risk types is given by the density $g(p)$, with support over $[\underline{p}, 1]$ for some $\underline{p} \in (0, 1)$. The expected return E is the same for all risk types. Under asymmetric information, the lender cannot discriminate between borrower risk types and therefore offers a pooling contract with gross interest rate r .

In this setting, [Ghatak \(1999\)](#) shows how the lender can harness joint-liability contracts in groups of two borrowers to mitigate credit rationing. Under this type of contract, a joint-liability payment $q \leq r$ is due in the asymmetric event where borrower i succeeds and partner j fails. [Ahlin and Townsend \(2007\)](#) extend this setting to allow for project returns that are positively correlated. This is implemented in the form of a constant $\bar{\epsilon}$ that adds probability mass to the symmetric events (where both borrowers succeed or fail) and subtracts it from the asymmetric events (where one group member fails and the other succeeds). In this model, the expected utility of borrower i forming a group with borrower j can be written as

$$u_{i,j} = E - rp_i - q[p_i(1 - p_j) - \bar{\epsilon}]. \quad (1)$$

Here, the expected utility is given by the expected project return E less the expected payable interest rp_i and expected joint-liability payment $q[p_i(1 - p_j) - \bar{\epsilon}]$. Because agents have no pledgeable collateral, borrower i only pays q in the asymmetric case where her project is successful and partner j defaults.

Agents face two decisions: (i) with whom and (ii) whether to take a loan. For the first decision, [Ghatak \(1999\)](#) shows that agents form groups that are homoge-

neous in risk type such that $p_i = p_j$. This follows from risk type complementary in Eqn 1, which exhibits a positive cross-partial derivative with respect to agents' risk types. For the second decision, agents take a loan when the expected utility $u_{i,j}$ exceeds that of the outside option \bar{u} . Because the cost of borrowing, i.e. the expected repayment, is strictly increasing in risk type, there is a marginal type \hat{p} that solves the participation equation

$$E - r\hat{p} - q[\hat{p}(1 - \hat{p}) - \bar{\epsilon}] = \bar{u} \quad (2)$$

with equality. Credit is rationed as borrowers with projects safer than \hat{p} do not find it profitable to borrow. Ahlin and Townsend (2007) argue here that increasing the project correlation mitigates credit rationing and thereby has a positive effect on the repayment to the bank. The intuition for this result is that higher $\bar{\epsilon}$ increases borrowers' utility by avoiding liability payments more often. This is because project correlation shifts probability mass from asymmetric to symmetric events. An increase in $\bar{\epsilon}$ therefore draws safer types into the market. This results in a new marginal type $\hat{p}' > \hat{p}$ and a safer borrower pool with types $p \in [\underline{p}, \hat{p}']$.

Developing the key trade-off

Contrary to the conclusions drawn in Ahlin and Townsend (2007), this safer borrower pool does not generally improve group repayment. To illustrate, note that after an increase in $\bar{\epsilon}$, the new marginal type \hat{p}' now has the same expected repayment ($E - \bar{u}$) as the previous marginal type \hat{p} . However, the previous marginal type and all others now have worse expected repayment (by the term $q \cdot d\bar{\epsilon}$) because the increase in correlation allows them to avoid liability payments more often. Proposition 3.1 provides conditions for project covariation to reduce repayment when the distribution of risk types is uniform.

Proposition 3.1. *Under a uniform distribution of risk types, the marginal effect of project covariation on expected repayment is strictly negative if either (i) the marginal type \hat{p} is smaller than $3/4$ or (ii) the joint-liability payment q does not exceed $3/5$ of the gross interest rate r .*

Proof: See Appendix A.

The intuition for the thresholds is that for correlation to improve repayment (i) the marginal types \hat{p} that are drawn into the market must be sufficiently safe to offset the negative effect of increased joint defaults and (ii) joint-liability payment q must be sufficiently high to lure the marginal types into the market in the first

place. Proposition 3.1 is limited to uniform distributions of risk types. Corollary 3.1 below shows that these thresholds are even higher for distributions with lower probability mass in the area of the marginal type.

Corollary 3.1. *The lower the density of the risk-type distribution $g(\hat{p})$ at the marginal risk type \hat{p} , the more an increase in project covariation will impair expected group repayment.*

Proof: See Appendix A.

The reasoning behind this corollary is that for an increase in project correlation to improve repayment, it must draw in considerably more safe types to offset the negative effect from borrowers avoiding joint-liability payments. For this to be the case, the distribution of types has to have considerable probability mass in the upper tail of the distribution.

Prediction for the context of the Bank for Agriculture and Agricultural Cooperatives

In the context of the Bank for Agriculture and Agricultural Cooperatives (BAAC), the model would predict a strictly negative repayment effect of correlation. The BAAC charges a fixed gross interest rate of 109% for small loans and joint-liability payments q are implemented in the form of a temporary increase in the payable interest rate. The maximum interest rate in the 1997 BAAC survey was 117%, which translates as a maximum joint-liability rate of $q = 8\%$ ($= 117\% - 109\%$). The ratio $q/r = 8\%/109\% \approx 0.07$ is well below the $3/5$ threshold. In addition, the actual distribution of types in the 2000 BAAC resurvey is Normal;² therefore, the predictions derived in Ahlin and Townsend (2007) cannot explain their empirical finding that repayment is better in markets with higher project correlation.

3.1.2 Ex-ante moral hazard: Direct effect

The Stiglitz (1990) model takes the homogeneous groups in Ghatak (1999) as given. The moral hazard problem relates to the following cooperative project choice after loan disbursement. Borrowers choose cooperatively between projects with different probabilities of success p_k with $k \in \{B, H\}$. Here B is the baseline project that was tied to the borrower in the previous subsection and H is the

²Shapiro-Wilk, Jarque-Bera and Kolmogorov-Smirnov tests of the risk-type variable (de-measured at the village level) cannot reject the null of Normality (N=292, p-values of 0.60, 0.65 and 0.81, respectively).

hazardous project with $p_H < p_B$. The hazardous project H has a higher *actual* outcome when successful, i.e. $y_H > y_B$, but a lower *expected* outcome, $p_H y_H < p_B y_B$. Information is asymmetric in that the lender does not observe which project is chosen but the group members do: Stiglitz assumes costless peer monitoring and enforcement. Group members make symmetric project choices that maximise their joint utility U_{kk} , resulting in individual project success probability

$$p = p_H \cdot 1[U_{BB} < U_{HH}] + p_B \cdot 1[U_{BB} \geq U_{HH}], \quad (3)$$

where $1[\cdot]$ is the Iverson bracket. In this context, the influence of project covariation $\bar{\epsilon}$ on the probability of repayment p depends on whether changes in $\bar{\epsilon}$ shift incentives towards the hazardous project. Using the project correlation structure introduced in the previous subsection, Ahlin and Townsend (2007) write the expected group pay-offs, given project choice $k \in \{B, H\}$, as

$$U_{kk} = U(y_k - r) \cdot [p_k^2 + \bar{\epsilon}] + U(y_k - r - q) \cdot [p_k(1 - p_k) - \bar{\epsilon}]. \quad (4)$$

Ahlin and Townsend (2007) now argue that the utility gain from avoiding joint-liability payment (of size $2q \cdot d\bar{\epsilon}$) due to an increase in $\bar{\epsilon}$ is comparatively higher for the baseline project, tilting incentives towards choosing the safer project. This is because (i) the baseline project has lower returns when successful and (ii) borrowers' utility is concave.

Developing the key trade-off

Again, the modelling in Ahlin and Townsend (2007) does not consider the negative effect that correlation has through borrowers avoiding joint-liability payments to the bank. The key trade-off is developed in Proposition 3.2 below.

Proposition 3.2. *The marginal effect of project covariation on repayment is strictly negative if either (i) borrowers are not extremely risk averse or (ii) the returns of the hazardous project are not substantially larger than those of the baseline project.*

Proof: See Appendix A.

The intuition for the negative repayment effect for risk-neutral borrowers is straightforward: with either (i) a linear utility function or (ii) $y_H \approx y_B$, the marginal increase in utility from higher project covariation is the same for both the baseline and the hazardous project, $\partial U_{BB}/\partial \bar{\epsilon} = \partial U_{HH}/\partial \bar{\epsilon} = 2q$. For (i), this is because

the slope of the utility function is constant. For (ii), this results from the gain in utility being evaluated at the same wealth level. Therefore, a change in $\bar{\epsilon}$ has no effect on project choice. However, it has a strictly negative effect of $-2q \cdot d\bar{\epsilon}$ from a diversification point of view because it reduces the probability that at least one borrower is successful.

3.2 Characterisation of stable matchings

This subsection extends the analysis of [Ghatak \(1999\)](#) by endogenising project correlation and allowing for a group size larger than two. Restricting this model to the empirical context, with two groups per market, I establish uniqueness conditions for stable matchings when utility is non-transferable. Equilibrium matching is shown to result in an endogeneity problem if borrowers' risk types are not fully captured by exogenous variables. I derive equilibrium conditions that impose simple inequalities on the latent group valuations and give traction to the empirical matching model that corrects for this bias.

3.2.1 Endogenous project correlation

The model used in the empirical application extends [Ghatak \(1999\)](#) to groups of size $n > 2$ and allows for project correlation that is determined endogenously. The latter is implemented by introducing three exposure types, A , B and N , which constitute the proportions θ_A , θ_B and θ_N of the agent population (as in [Ahlin, 2009](#)). N -types are not affected by external shocks. For A - and B -types, the independent shocks A and B equiprobably add or subtract, respectively, a shock term γ from the project success probability. Extending the model in Eqn 1 in this way, borrower i 's utility from taking a loan with group G can be written as

$$u_{i,G} = E - rp_i - qp_i \sum_{j \in G \setminus i} (1 - p_j) + q\epsilon \sum_{s \in \{A,B\}} 1[i \in s] \cdot (n_s^G - 1), \quad (5)$$

where $1[\cdot]$ is the Iverson bracket, which is 1 if borrower i is of exposure type $s \in \{A, B\}$ and 0 otherwise, n_s^G is the number of borrowers of exposure type s in group G and the constant $\epsilon := \gamma^2$ gives the intensity of the projects' exposure to shocks.

3.2.2 Assumptions

The analysis below makes four assumptions. First, reflecting the nature of the Thai group-lending data (see Section 5), the analysis is restricted to two groups per market. Second, I treat the distribution of risk types $p \in [p, 1)$ and exposure types $s \in \{N, A, B\}$ as independent. The spineplot in Figure 2a plots these two variables against each other based on the 2000 BAAC survey used in the empirical analysis. The exposure types on the vertical axis are categorised based on which of the previous two years was worse for the borrower economically. The data exhibits no systematic relationship between risk type p on the horizontal axis and exposure type s on the vertical axis. I therefore assume

$$p \perp s. \tag{H1}$$

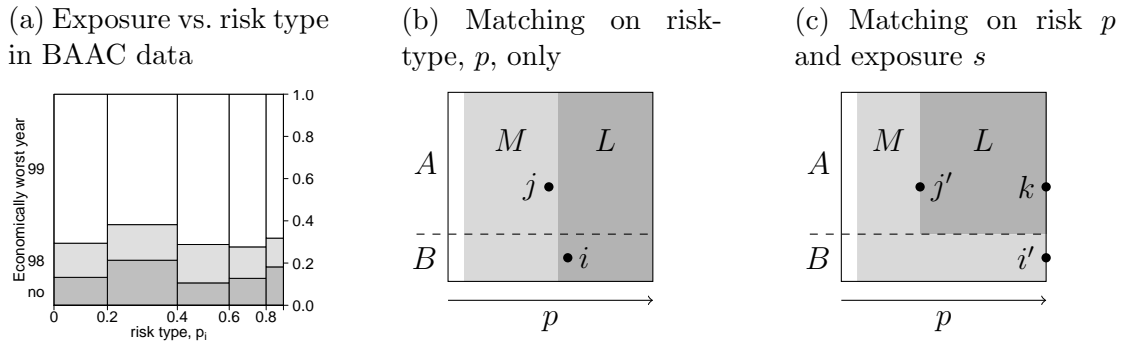
Third, I assume that utility is not linearly transferable between borrowers. In models with non-transferable utility, agents always prefer matches with higher valuations. That is, agents cannot negotiate binding contracts to compensate for committing to match with less attractive partners. While this assumption is less common in the microfinance literature, there is no empirical evidence for the existence of such transfers. Furthermore, the assumption of non-transferable utility is particularly well placed in the context of group lending, where (i) ex-ante transfers are not possible due to limited initial wealth and (ii) [Holmström and Tirole \(1997\)](#) show that incentives are muted if a borrower initially pledges too much of her future income.

Lastly, for the likelihood of the empirical model to be well defined, the observed equilibrium in the data must be a unique stable matching ([Bresnahan and Reiss, 1991](#)). The existence and uniqueness of an equilibrium can be guaranteed by imposing suitable restrictions on agents' preferences. This is common practice in the empirical analysis of matching markets (see the literature review). I build on results in [Pycia \(2012\)](#), who shows that pairwise-aligned preferences are both necessary and sufficient for the existence of a unique equilibrium matching. Pairwise-aligned preferences imply that any two borrowers that belong to the same two groups prefer the same group over the other. That is, for an equilibrium group ABC , aligned preferences would imply that borrowers A and B agree on the relative ranks of C and D , i.e. $ABC \succsim_A ABD \Leftrightarrow ABC \succsim_B ABD$, where \succsim_i represents agent i 's preference relation over groups that contain i .

3.2.3 Equilibrium characterisation

Figure 2b illustrates pairwise-aligned preferences for two groups, M and L . Here, matching is on risk type p only. This is equivalent to assuming that the measure of risk exposure intensity ϵ is 0. In this case, preferences are aligned in that the groups are strictly rank-ordered by risk type (because of risk-type complementarity). In the following, I refer to the group with the highest risk-ordering of types as the *dominant* group L (dark shading) and the other group as the *residual* group M (light shading).

Figure 2: Matching on risk type (horizontal axis) and exposure type (vertical axis) in two-group markets.



For $\epsilon > 0$, sorting takes place along two dimensions, where preferences are

- *aligned* in risk type p , in that borrowers always prefer a safer partner (irrespective of their own type), but also
- *assortative* in exposure type s , in that borrowers only value partners of their own type.

In two-group markets, the existence of a unique equilibrium is guaranteed if preferences are aligned in the dominant group L . This is because no member of this group will find it attractive to switch to residual group M and therefore the matching is stable. Proposition 3.3 derives the necessary conditions.

Proposition 3.3. *In two-group markets, preferences are aligned in the dominant group L if either (i) ϵ is zero – or, equivalently, all agents are of the same exposure type – or (ii) ϵ and the proportion of the leading exposure type A are sufficiently large.*

Proof: See Appendix A.

The conditions in Proposition 3.3 are reasonable in an agricultural context where rice farmers are the dominant group (as in the Townsend Thai project, see Section 5) and their projects are arguably subject to intense common shocks. The intuition here is that in agricultural lending, there is one dominant exposure type: exposure to weather shocks. If these shocks are sufficiently strong, then matching is first on exposure type and then aligned in risk type within exposure type (as Figure 2c illustrates). In villages with two borrower groups, this results in a dominant group, which is composed of the ‘weather shock’ exposure types with the safest projects (here, group L). This group is stable (or in equilibrium) when utility is non-transferable, because no borrower would prefer to match with a different exposure type or a higher risk type. The remaining borrowers form a residual group, which is composed of ‘weather shock’ exposure types with riskier projects and those with other exposure types (here, Group M). At the same time, this equilibrium matching (i) creates an endogeneity problem that results in sorting bias and (ii) provides an elegant solution to the problem. Both results are discussed in turn.

Sorting bias

Aligned preferences result in the maximisation of the dominant group’s valuation. The valuation V_G of group G is the sum over all group members’ utilities from matching with this group, i.e. the sum over the interaction terms in Eqn 5.

$$V_G = -q \sum_{i \in G} \sum_{j \in G \setminus i} [p_i(1 - p_j) + p_j(1 - p_i)] + q\epsilon \sum_{s \in \{A, B\}} n_s^G (n_s^G - 1) \quad (6)$$

In particular, group valuation V_G does not contain borrower i ’s expected return E and interest payment rp_i , because these realisations are independent of group members. Group valuation in Eqn 6 is increasing in risk type (for $p > 0.5$), exposure intensity ϵ and the coincidence of same exposure types. Equilibrium matching, as illustrated in Figure 2c, results in a positive correlation between the groups’ risk type (first term of Eqn 6) and project covariation (second term). Figure 2c shows the equilibrium matching, where the dominant group is homogeneous in exposure type (all A types). We can see that the group with higher project covariation (group L) has safer risk types on average. Corollary 3.2 states this formally.

Corollary 3.2. *In two-group markets with project correlation, equilibrium matching exhibits positive correlation between the groups’ average risk type and project*

covariation.

Proof: See Appendix A.

In the resulting matching, the dominant group has, on average, projects that are both safer and more highly correlated than the residual group. Now, if risk type is partially unobservable and captured in the error term, then the coefficient pertaining to project correlation will be biased. I refer to this as sorting bias throughout the paper.

Equilibrium characterisation

The equilibrium conditions can be expressed as simple inequalities that impose lower and upper bounds on the match valuations of the observed and unobserved (or counterfactual) matches. I impose these bounds in the empirical matching model in Section 4 to guarantee that a unique equilibrium is estimated. Proposition 3.4 summarises the stability conditions based on bounds \overline{V}_G and \underline{V}_G , derived in Appendix A. The proof is for aligned preferences in the general case with arbitrary group and market size. The conditions for observed equilibrium groups $G \in \mu$ and unobserved non-equilibrium groups $G \notin \mu$ are equivalent, but they impose different bounds on the latent valuation variables that guarantee the estimation of the unique market equilibrium.

Proposition 3.4. *The matching μ is stable iff $V_G < \overline{V}_G \quad \forall G \notin \mu$. Equivalently, the matching μ is stable iff $V_G > \underline{V}_G \quad \forall G \in \mu$.*

Proof: See Appendix A.

The upper bounds \overline{V}_G have a natural economic interpretation; they are the maximum of the opportunity costs for group G 's members of leaving their respective equilibrium groups and joining non-equilibrium group G . Similarly, the lower bounds \underline{V}_G give the maximum of the opportunity costs of group G 's members maintaining their equilibrium match G .

4 Empirical strategy

This section outlines the empirical strategy used to identify the direct and participation effects separately. I describe what is being tested in the following and how these tests relate back to the theoretical models.

4.1 Direct effect

Subsection 4.1.1 develops a structural empirical model to estimate the direct repayment effect of project correlation net of sorting bias. The estimation strategy replicates the following ideal experiment with standard cross-sectional survey data.

Ideal Experiment 1: Direct effect net of sorting bias

1. Announce in each village that loan applicants will be assigned to groups randomly and make applicants sign up to a waiting list.
2. For half of the villages (chosen at random), surprise applicants by allowing groups to form endogenously. For the other half, assign groups randomly.
3. Obtain the parameter estimates of randomly and endogenously formed groups. Call the first estimates the *direct* effect of project covariation and the difference between the two groups the *bias from sorting*.

4.1.1 Estimation strategy

Technically, the equilibrium groups constitute a self-selected sample. The selection problem differs substantially from the classical Heckman (1979) two-stage correction. Here, the first-stage selection mechanism that determines which borrower groups are observed (and which are not) is a one-sided matching game and not a simple discrete choice, as in the Heckman model. A discrete choice model assumes that an observed match reveals group partners' preferences concerning each other. An observed matching, however, is the outcome of complex interactions and conflicts of interest between agents. In particular, borrowers can only choose from the set of partners who would be willing to form a match with them, but we do not observe their relevant choice sets. This makes direct inference based on a discrete choice model impossible, even if it accounts for social interactions such as the models in Brock and Durlauf (2007) and Ciliberto and Tamer (2009).

The empirical strategy, therefore, is to simultaneously estimate the outcome equation of repayment performance with the matching game. The matching game is given by the following match equation

$$V_G = W_G\alpha + \eta_G. \quad (7)$$

There are $|\Omega|$ equations, where Ω is the set of feasible groups in the market.³

³The set of feasible groups in two-group markets with group size n comprises all $\binom{2n}{n}$ possible k -for- k borrower swaps for $k \in \{1, \dots, n-1\}$ across the two groups.

$V \in \mathbb{R}^{|\Omega|}$ is a vector of latents and $W \in \mathbb{R}^{|\Omega| \times k}$, a matrix of k characteristics for all feasible groups. $\alpha \in \mathbb{R}^k$ is a parameter vector, and $\eta \in \mathbb{R}^{|\Omega|}$ is a vector of random errors. A group – and therefore its repayment outcome Y_G – is observed if it is part of the equilibrium matching μ , i.e. its group valuation is in the set of valuations Γ_μ that satisfy the equilibrium condition.⁴ This set of valuations is the link between the structural empirical model and the equilibrium characterisations derived in Proposition 3.4, Subsection 3.2. With $V \in \mathbb{R}^{|\Omega|}$, the vector of all valuations in the market, the equilibrium condition can be written as a collection of inequalities that give upper and lower bounds on the match valuations

$$V \in \Gamma_\mu \Leftrightarrow [V_G < \overline{V}_G \ \forall G \notin \mu] \Leftrightarrow [V_G > \underline{V}_G \ \forall G \in \mu]. \quad (8)$$

For the outcome equation, the binary dependent variable is given as $Y_G = 1[Y_G^* > 0]$, where the latent group outcome variable Y_G^* is

$$Y_G^* = X_G \beta + \varepsilon_G, \quad (9)$$

with $\varepsilon_G := \delta \eta_G + \zeta_G$, where ζ_G is a random error. This specification allows for a linear relationship between the error terms in the selection and outcome equations with covariance δ . The design matrices $X \in \mathbb{R}^{|\mu|}$ and $W \in \mathbb{R}^{|\Omega|}$ do not necessarily contain distinct explanatory variables.

Distribution of error terms

The joint distribution of ε_G and η_G is assumed bivariate normal with mean zero and constant covariance δ .

$$\begin{pmatrix} \varepsilon_G \\ \eta_G \end{pmatrix} \sim N \left(0, \begin{bmatrix} \sigma_\xi^2 + \delta^2 & \delta \\ \delta & 1 \end{bmatrix} \right) \quad (10)$$

Here, the variance of the error term of the outcome equation σ_ε^2 is $\text{var}(\delta \eta + \xi) = \delta^2 + \sigma_\xi^2$. To normalise the parameter scale, the variance of η and ζ is set to 1, which simplifies σ_ε^2 to $1 + \delta^2$ in the estimation. If the covariance δ were zero, the marginal distributions of ε_G and η_G would be independent and the selection problem would vanish.

⁴The Heckman (1979) model is a special case where the set of feasible valuations is $\Gamma = [0, +\infty)$.

Identification

Interaction in the market makes estimation computationally involved but also overcomes the identification problem. Identification requires exogenous variation. In this model, this is provided for every group by the characteristics of agents who are in the same market but not in the same group. To illustrate, take a market with four agents A , B , C and D . The characteristics in the outcome equation of group AB are simply $X = (X_{AB})$. The characteristics in the matching equation are $W = (X_{AB}, X_{CD}, X_{AC}, X_{AD}, X_{BC}, X_{BD})$, and the independent elements of W are then $W' = (X_{CD}, X_{AC}, X_{AD}, X_{BC}, X_{BD})$. The identifying assumption is thus that the characteristics of agents outside the match (those comprised in W') are exogenous. Put differently, the identifying exclusion restriction is that the characteristics of all agents in the market affect who matches with whom, but the outcome of an equilibrium group is determined exclusively by its own members. Note that other agents' characteristics are not used as instruments in a traditional sense. Rather than entering the selection equation directly, they pose restrictions on the match valuations by determining the bounds in the estimation.

Estimation

In the estimation, I follow [Sørensen \(2007\)](#), who uses Bayesian inference with a Gibbs sampling algorithm that performs Markov Chain Monte Carlo (MCMC) simulations from truncated normal distributions. The latent outcome and valuation variables Y^* and V are treated as nuisance parameters and sampled from truncated Normal distributions that enforce sufficient conditions for the draws to come from the equilibrium of the group formation game. For the posterior distributions, see [Appendix B](#). The conjugate prior distributions of parameters α , β and δ are Normal and denoted by $N(\bar{\alpha}, \Sigma_\alpha)$, $N(\bar{\beta}, \Sigma_\beta)$ and $N(\bar{\delta}, \sigma_\delta^2)$, respectively. In the estimation, the prior distributions of α and β have mean zero and variance-covariance matrix $\Sigma_\beta = (\frac{1}{|\mu|} X'X)^{-1}$ and $\Sigma_\alpha = (\frac{1}{|\Omega|} W'W)^{-1}$, respectively. This is the widely used g-prior ([Zellner, 1986](#)). For δ , the prior distribution has mean zero and variance 10. For this parameter, the prior variance is at least 40 times larger than the posterior variance in all estimated models. This confirms that the prior is fairly uninformative.

4.1.2 Testable effects and links to theory

By linking the structural empirical model to the variables defined in the theory, the empirical specification of the matching and outcome equations can be written

as

$$V_G = -q \sum_{i \in G} \sum_{j \in G \setminus i} [p_i + p_j - 2p_i p_j] + q\epsilon \sum_{s \in \{A, B\}} n_s^G (n_s^G - 1) + \eta_G \quad (11)$$

$$Y_G^* = r \sum_{i \in G} p_i + q \sum_{i \in G} \sum_{j \in G \setminus i} [p_i + p_j - 2p_i p_j] - q\epsilon \sum_{s \in \{A, B\}} n_s^G (n_s^G - 1) + \delta \eta_G + \zeta_G. \quad (12)$$

The matching equation is the empirical equivalent of Eqn 6. Eqn 12 gives the expected repayment Y_G^* of group G . In words, the expected repayment equals the expected interest payment plus the expected liability payment (if projects are independent) and minus the liability payment that the group avoids due to correlated returns. The final term $\delta \eta_G$ controls for unobservable group characteristics through the error term of the matching equation η_G . The error term ζ_G captures realised individual or aggregate shocks such as health or market demand effects.

For the parameters, the gross interest rate r is known to be fixed at 1.09 in the BAAC lending programme and is therefore fixed at this level here. The parameters q , $q\epsilon$ and δ are estimated in the model. The expected signs of the parameters are as given in Eqns 11 and 12. Of particular interest is the sign of $q\epsilon$, which pertains to the project correlation variable in the outcome equation. From a diversification point of view, project correlation has a strictly negative effect. However, this effect can be (i) outweighed by a positive effect from mitigating moral hazard (see Proposition 3.2) or (ii) confounded by a positive sorting bias from endogenous group formation (see Corollary 3.2.). Controlling for unobservable group valuation η_G allows me to estimate the direct repayment effect net of sorting bias. The extent and sign of the sorting bias are captured by parameter δ .

4.2 Participation effect

In a second step, I test for the participation effect of restricting matching on risk exposure. This effect is estimated in agent-based simulations using the coefficient estimates from the structural model as parameters. Varying the matching process but keeping the model fixed to predict the repayment outcome allows me to separate the participation effect from the direct effect. The agent-based simulations can be thought of as replicating the following ideal experiment.

Ideal Experiment 2: Participation effect

1. Randomly assign villages to one of two regimes. Dependent on the regime, have groups apply under either (i) matching on risk type only – i.e. groups must be balanced in exposure type – or (ii) matching on both risk and exposure type.
2. For all villages, surprise loan applicants by disbursing individual-liability loans instead of joint-liability loans.
3. Compare the average repayment rates under the two regimes. Call the difference in repayment the participation effect of matching on risk exposure.

4.2.1 Estimation strategy

To estimate the size of the participation effect, I work with the full sample of borrowers in the 2000 BAAC data and run agent-based simulations to see how many and what sorts of groups will borrow at the current contract terms under different matching regimes. The characteristics of these self-selected groups are then used to predict the expected repayment using the parameter estimates from Eqn 12. The agent-based simulation follows the protocol below.

1. Obtain the equilibrium groups in the 29 two-group markets for different matching regimes: (i) matching on risk type only and (ii) matching on both risk and exposure type. Equilibrium groups are determined using the group valuation in Eqn 11 (with η_G set to zero) and equilibrium conditions derived in Proposition 3.4.
2. Calculate borrower i 's expected pay-off $\tilde{u}_{i,G}$ from taking a loan with equilibrium group G based on the empirical specification of Eqn 5 as follows

$$\tilde{u}_{i,G} = E_i + \bar{l}_t \left[1 - rp_i - \hat{q}p_i \sum_{j \in G \setminus i} (1 - p_j) + \hat{q}\epsilon \sum_{s \in \{A,B\}} 1[i \in s] \cdot (n_s^G - 1) \right], \quad (13)$$

where \bar{l}_t is the median loan size in market t and \hat{q} and $\hat{q}\epsilon$ are the parameter estimates from Eqn 11.

3. Evaluate each borrower's participation condition $\tilde{u}_{i,G} > \bar{u}_t$, where \bar{u}_t is the value of a borrower's outside option, measured as the median wage rate for agricultural labour in that market.

4. Exclude a group from the sample if the participation condition is satisfied for fewer members than the minimum group size in the market.
5. For the remaining groups, predict the expected repayment using Eqn 12 and the parameter estimates from the structural empirical model.

4.2.2 Testable effects and links to theory

Similar to that for the direct effect, the test for the participation effect is positioned between two opposing predictions. On the one hand, the Ghatak (1999) adverse selection model results in a negative repayment effect. On the other hand, Katzur and Lensink (2012) show that the perfect information outcome can be achieved – in the Ghatak (2000) model with a binary distribution of risk types – if project covariation is sufficiently high for safe groups compared to risky groups. While it is not clear that the result in Katzur and Lensink (2012) carries over to our context, it still merits consideration when interpreting the results.

5 Empirical results

The empirical strategy in Section 4 is applied to data from the Townsend Thai project. The analysis here uses data from both the 1997 baseline survey and a smaller resurvey conducted in 2000. Replication code and datasets are available in R package `matchingMarkets` (Klein, 2015b) and the corresponding Vignette (Klein, 2015a).

5.1 Data

The survey project is a panel that focuses on villages in four provinces (*changwat*) of Thailand: two in the North-east region and two in the Central region. The baseline data used in the Ahlin and Townsend (2007) paper was collected in 1997. For this study, 12 subdistricts (*tambons*) were selected at random within each of the four provinces. Within each *tambon*, four villages were selected at random. This resulted in a sample of 192 villages, in which two survey instruments were applied. In the initial household survey (Townsend, 1997b), 15 households in each village were selected at random, yielding a total sample of $192 \times 15 = 2,880$ households. The second survey instrument was the initial Bank for Agriculture and Agricultural Cooperatives (BAAC) survey (Townsend, 1997a) or BAAC 1997. The BAAC is a government-owned development bank and the largest lender to

this population. In the BAAC 1997 survey, for every village as many borrower groups as possible were identified and a maximum of two groups were randomly selected for interviews. In total, 262 BAAC groups were identified and their group leaders interviewed.

For the main part of the analysis, I use data from a smaller resurvey that was conducted in 2000 and comprises variables that were specifically designed to test the theory in Ghatak (1999). In the resurvey, for each of the four original provinces, four *tambons* were selected randomly from the 12 *tambons* in the baseline survey. This resurvey again consisted of two instruments: a household resurvey (Townsend, 2000b) and a BAAC resurvey (Townsend, 2000a), referred to as BAAC 2000 in the following. BAAC 2000 consists of a group-leader survey, in which the heads of BAAC groups were interviewed, as well as a group survey, in which up to five group members were interviewed. The final sample of the BAAC 2000 used for analysis comprises the characteristics of 68 lending groups.

Table 2: Summary of group-level variables.*

Variable	Description	mean (sd)
<i>Dependent variable</i>		
- repayment_outcome ^{a)}	BAAC <i>never</i> raised interest rates as a penalty for late repayment	0.46 (0.50)
<i>Exposure</i>		
- ln(group_age) ^{b)}	Log of number of years group had existed	2.33 (0.55)
<i>Risk type</i>		
- success_prob $p_i^c)$	Group members' project success prob.	0.41 (0.15)
- success_prob_int $p_i p_j^c)$	Two-way interactions of success prob.	0.24 (0.05)
<i>Project covariation</i>		
- worst_year $wst^c)$	Measure of coincidence of economically bad years across group members	0.57 (0.37)
<i>Contract terms</i>		
- interest_rate	Gross interest rate is fixed at 109% for loans below 60,000 Thai baht	1.09 (–)
- loan_size ^{a)}	Average loan size borrowed by the group (thousand Thai baht, currency value in 2000)	17.12 (10.87)

^{a)} from 2000 BAAC group-leader survey

^{b)} random regression imputation based on 1997 and 2000 BAAC surveys

^{c)} from 2000 BAAC group survey

5.2 Variables

The variables used in the empirical analysis are directly related to the extension of Ghatak's (1999) theoretical model of borrower group formation in Subsection 3.2. The average *risk type* and *project covariation* are measured as below, and the remaining variables are summarised in Table 2.

Risk type: Group members were asked for their expected income for the following year, which is denoted as E_i . They were also asked for their expected income if the following year was a good year H_i or a bad year L_i . The measure $p_i = \frac{E_i - L_i}{H_i - L_i}$ serves as a proxy for borrower i 's probability of success, using the property that $p_i H_i + (1 - p_i) L_i = E_i$.

Project covariation: A group's project covariation is proxied by the variable *worst_year*, which is a vector indicating which of the previous two years was worse for a borrower economically. The group-level variable gives the average coincidence of worst years based on all possible borrower-by-borrower comparisons. This measure establishes a direct link with the different exposure types in Ahlin (2009) in that each year can then be interpreted as exposing agents to a different shock. The measure of project covariation then gives the probability that two randomly drawn group members have the same exposure type.

5.3 Direct effect

The first Probit model in Table 3 gives the marginal effect of project covariation on repayment. The dependent variable is 1 if there were no arrears during the group's lifetime and 0 otherwise. To compare the riskiness of groups with different ages and, therefore, different exposure to risk, I control for the natural logarithm of group age. I also add village-level fixed effects to control for between-village heterogeneity. The resulting positive coefficient suggests that a high level of project covariation is associated with less arrears. This replicates the surprising result in Ahlin and Townsend (2007) using data from the 2000 resurvey. This positive repayment effect can be explained by either correlation mitigating moral hazard for extremely risk-averse borrowers (see the Stiglitz model, Proposition 3.2) or the endogenous matching that biases $\hat{\beta}_{wst}$ upwards because it picks up the effect of the omitted risk-type variable (see Corollary 3.2). To explore this bias from sorting, the second Probit model controls for contract terms and the positive repayment effect of risk type. This control mitigates the sorting bias for $\hat{\beta}_{wst}$ and results in a switch in sign, which is consistent with the negative effect from anti-diversification, as predicted in the Stiglitz model for moderately risk-averse agents (see Eqn 12).

Table 3: Probit and structural models with market fixed effects

<i>S.E. in parentheses; one-sided significance at 0.1, 1, 5, 10% denoted by ***, **, *, and .</i>			
	Probit model (1)	Probit model (2)	Structural
Outcome equation			
<i>Dependent variable: repayment_outcome^{a)} = 1 if the BAAC has never raised interest as a penalty for late repayment; 0 otherwise.</i>			
<i>Risk type</i>			
- success_prob p_i	–	+1	+1
- success_prob_int $p_i p_j$	–	0.238 (1.606)	1.571 (1.813)
<i>Project covariation</i>			
- same_worst_year $wst^a)$	0.170 (0.289)	-0.015 (0.219)	-0.586 (0.243) **
<i>Contract terms</i>			
- loan_size	–	0.263 (0.421)	0.970 (0.362) **
- loan_size_sqrd	–	-0.050 (0.088)	-0.187 (0.080) *
<i>Exposure</i>			
- ln(group_age)	-0.040 (0.054)	-0.116 (0.161)	-0.395 (0.109)***
<i>Village-level controls</i>			
	YES	YES	YES
Observations	68	68	68
Matching equation			
<i>Dependent variable: group observability indicator = 1 if group is observed; 0 otherwise.</i>			
<i>Risk type</i>			
- success_prob_int $p_i p_j$	–	–	-0.778 (0.992)
<i>Project covariation</i>			
- same_worst_year wst	–	–	0.356 (0.119) **
<i>Controls</i>			
	–	–	YES
Observations ^{b)}	–	–	5,342
Variance			
Covariance δ	–	–	0.512 (0.127)***

^{a)} [Karlson et al. \(2012\)](#) one-sided test for difference of *Probit(2)* and *Structural*, p-value 0.048.

^{b)} 5,284 counterfactual groups and 58 factual groups.

5.3.1 Matching on observables

The above switch in sign implies a positive correlation between risk type and exposure type, which results from endogenous matching on both covariates as derived in Corollary 3.2. To confirm that this is the mechanism at work, the matching on observables is tested in the matching equation of the structural model in Table 3. In this equation, the independent variables are constructed from individual borrowers' characteristics for 58 factual (or equilibrium) groups and 5,284 counterfactual (or non-equilibrium) groups in all 29 two-group villages. The dependent variable is 1 for the 58 equilibrium groups and 0 otherwise. The latent group valuations are simulated for equilibrium and non-equilibrium groups using

the Gibbs sampler presented in Subsection 4.1.1. Turning to the results, the signs of the marginal effects⁵ are consistent with the predictions from the theory in Eqn 11.⁶ The negative sign on the risk-type variable means that borrowers value group members with safer projects. (Note: the negative sign on the coefficient results in a positive cross-partial derivative with respect to agents' risk types in Eqn 11.) While this effect is non-significant, the positive sign on the exposure-type variable is significant at the 1%-level and indicates that borrowers value peers of the same exposure type. This finding is in line with the matching mechanism derived in Corollary 3.2. This is interesting in that it suggests that exposure type may play an even more significant role in group formation than risk type, which has been the primary focus of the microfinance literature to date.

5.3.2 Matching on unobservables

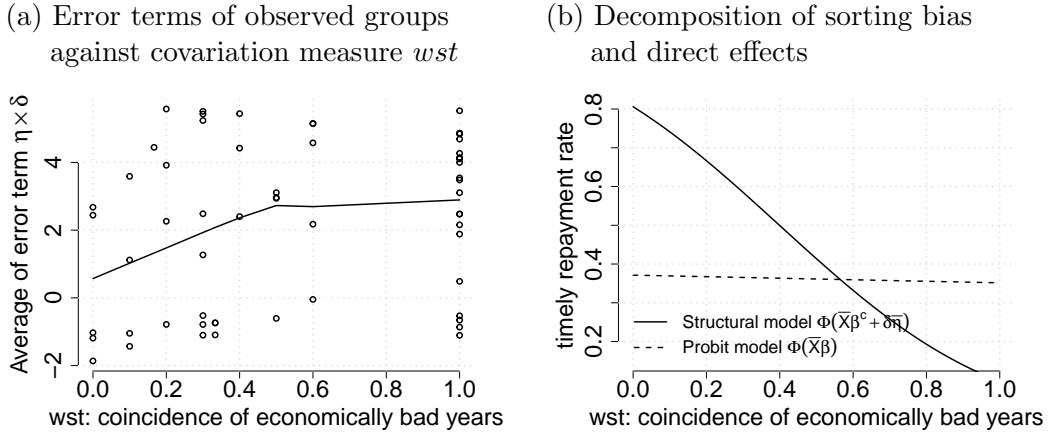
If matching is also on unobservables that affect group repayment – such as local information on risk types – then $\hat{\beta}_{wst}$ is still biased upwards in the second Probit model. To correct for this bias, the structural model in Table 3 estimates the matching and outcome equations jointly and allows local information to enter the outcome equation in form of the error term η of the matching equation. The variance section in Table 3 shows considerable matching on unobservables: the covariance between the error terms of the matching and outcome equations is $\hat{\delta} = 0.512$, which is equivalent to a correlation of $+0.41$ ($= \frac{\sigma_{\varepsilon, \eta}}{\sigma_{\varepsilon} \sigma_{\eta}} = \frac{0.512}{(1+0.512^2) \cdot 1}$). A direct comparison between the second Probit model and the sorting-corrected structural model yields an upwards bias in the Probit model of $+0.57$ ($= -0.015 - [-0.586]$) for $\hat{\beta}_{wst}$ that is significant at the 5%-level. This bias results from the positive correlation of project covariation and unobservables η in the outcome equation (see Figure 3a). In the case of group lending, this means that groups with higher project covariation also have better unobserved characteristics. In the structural model, the error term η in the matching equation enters the outcome equation as $\delta\bar{\eta} > 0$. The omission of this sorting-correction term in the Probit regression leads

⁵The marginal effects for the selection equation are obtained as $\frac{\partial P}{\partial W} = \phi(0)\alpha/\sqrt{2}$, with the probability P that group G has a higher valuation than group G' equal to $Pr(W_G\alpha + \eta_G > W_{G'}\alpha + \eta_{G'}) = \Phi((W_G - W_{G'})\alpha/\sigma_{\eta_G - \eta_{G'}}) = \Phi((W_G - W_{G'})\alpha/\sqrt{2})$. The standard error of the marginal effect is given by $\phi(0)\sigma_{\alpha}/\sqrt{2}$. To see this, consider a linear transformation of $X \sim N(\mu, \sigma)$ as $Y = aX$. It then follows that $Y \sim N(a\mu, a\sigma)$.

⁶Note that coefficient magnitudes on risk type and correlatedness need not be the same in both equations of the structural model. This is mainly for two reasons: first, the response of the outcome equation is the probability of timely repayment rather than the expected repayment; and second, the parameters in the outcome equation are based on the adverse selection model and do not reflect the moral hazard effects through which correlated projects can tilt incentives towards safer projects.

to a positive correlation $cor(wst, \varepsilon)$ because ε is proportional to $\delta\eta$ (i.e. $\varepsilon = \delta\eta + \xi$, where ξ is a random error). Matching on both observables (wst) and unobservables (η) thus explains the sorting bias in the second Probit model.

Figure 3: Matching on unobservables. Relative magnitudes of sorting bias and the direct effect of project covariation on repayment outcomes.



5.3.3 Decomposition of sorting bias and direct effect

Figure 3b illustrates the decomposition of sorting bias and direct effects. The decomposition is done by comparing the estimated regression lines for the first Probit model with the outcome equation of the structural model. The models are evaluated conditional on the value of wst on the horizontal axis with all other variables at their means. The solid regression line of the structural model gives the expected repayment – conditional on wst – when all borrowers are randomly assigned to groups. This is because the estimates are conditional on all feasible groups (observed and unobserved) in the market. The dashed Probit regression line depicts the estimates for observed groups only and therefore also captures the sorting bias. To emphasise, if borrowers were assigned at random, as in Ideal Experiment 1 in Section 4, the two lines would overlap perfectly.

In Figure 3b, we see that allowing groups to match endogenously (dashed line) results in more timely repayment for groups with higher project covariation. This is the result in Ahlin and Townsend (2007). However, it does not imply a causal relationship. To quantify this effect, note that an increase in project covariation by one standard deviation at the sample mean results in an expected improvement in the probability of timely repayments of +6.3 percentage points ($= 0.170 \cdot 0.37 = \hat{\beta}_{wst}^{probit} \cdot \hat{\sigma}_{wst}$). This improvement follows from two opposing

effects. First, from the structural model we find a significantly negative *direct effect* of -22 percentage points ($-0.586 \cdot 0.37 = \hat{\beta}_{wst}^{str} \cdot \hat{\sigma}_{wst}$) because the bank loses joint-liability payments when projects fail simultaneously. This is consistent with the revised predictions from the moral hazard model of [Stiglitz \(1990\)](#) when borrowers are not extremely risk averse. Second, from the difference between the Probit and structural models we find an even larger but positive *sorting bias* of $+28$ percentage points ($(\hat{\beta}_{wst}^{probit} - \hat{\beta}_{wst}^{str}) \cdot \hat{\sigma}_{wst}$). This is because the highly correlated groups have unobservables that make them $+28$ percentage points safer.

5.3.4 Robustness of the results

In this subsection, I examine whether my primary result – the decomposition into a negative direct effect and a positive sorting bias – is robust to various empirical issues.

I first examine a potential reverse causation problem, in that all group members may report their worst year as that in which their group faced repayment problems. This would provide an alternative explanation as to why groups with correlated returns have worse repayment outcomes. To rule out this explanation, first note that when borrowers were asked why they perceived one year as worse than another, only five out of a total of 390 borrowers gave the reason ‘unable to repay debt’ in their response. In addition, repayment was surveyed retrospectively over the full lifetime of groups. The average group age was 11 years, but project correlation is calculated based on just two years.

A second concern is survivorship bias. Groups with safer types are more likely to survive, particularly when returns are highly correlated. This ‘survival of the safest’ would result in groups with more correlated returns being safer and provide an alternative to my endogenous matching explanation. To disprove this explanation, it is enough to show that older groups are not safer on average. In fact, the correlation between risk type and group age is negative, at -0.065 , p -value 0.599 , meaning that survivorship bias is not an issue.

Finally, the equilibrium conditions are derived based on the assumption that the matching data represent the complete market. In the paper, the model is estimated using a random sample of five borrowers from groups with 11 borrowers on average. This is a shortcoming in the empirical analysis. However, [Klein \(2015a\)](#) presents Monte Carlo evidence of the robustness of the estimator in small samples, which confirms that the resulting attenuation bias even underrates the sorting bias that this paper corrects for.

5.4 Participation effect

For the direct effect, the empirical model does not allow for an outside option leaving in or excluding some potential borrower groups. The participation effect tests whether allowing for matching on exposure type can draw sufficiently safe types – that would not have taken a loan otherwise – into the market in order to offset the negative effect from avoiding liability payments. This is an indirect test of the model extension of Ghatak (1999) in Subsection 3.1, which predicts a negative repayment effect, against Katzur and Lensink (2012), who show that group lending can achieve the perfect information outcome if project covariation is sufficiently higher for safe groups compared to risky groups.

Table 4: Agent-based simulation of expected repayment under different matching regimes.

Simulation based on 250 individuals in 29 two-group markets.		
<i>Matching process:</i>	matching on p only (1)	matching on p and s (2)
<i>Participation</i>		
1. No. of borrowers	165	164
2. No. of groups	17	17
<i>Group characteristics</i>		
3. \bar{p}	0.740	0.723
4. \bar{wst}	0.571	0.639
<i>Predicted repayment</i>		
5. \hat{Y}	0.428	0.390
6. 80% CI ^{a)}	(0.904, 0.048)	(0.905, 0.031)

^{a)} Confidence intervals based on endpoint transformation.

Table 4 presents the results of the simulations for (1) matching on risk type only versus (2) matching on both risk type and exposure type. The first row gives the number of borrowers whose utility from taking a loan with their equilibrium groups exceeds the outside option of wage labour. Contrary to the predictions of the theories, matching on exposure type (anti-diversification) does not draw more borrowers into the programme (164) than matching on risk type alone (165). This result carries through from the individual to the group level: the restriction on minimum group size makes borrowing infeasible for groups where the participation condition is satisfied for fewer borrowers than the minimum group size. The number of remaining groups is given in the second row.

While anti-diversification does not draw in more borrowing groups (17 vs 17), these groups are riskier (\bar{p} in Row 3) and have considerably higher project correlation ($\bar{w}st$ in Row 4). The predicted probability of timely repayment of 0.390 for these groups is consequently lower than when matching on risk type only (0.428). This is because under high project correlation, the bank receives fewer joint-liability payments, consistent with the model predictions from [Ghatak \(1999\)](#). The effect is statistically insignificant but of economic importance: anti-diversification results in a 10 percent increase in timely repayment. The results further suggest that the predictions from [Katzur and Lensink \(2012\)](#) are not applicable in this context and that lenders would benefit from preventing the grouping together of borrowers exposed to similar income shocks.

5.5 Implications for market design

The results concerning the direct and participation effects imply that banks should prevent the matching of borrowers who are exposed to similar income shocks. A policy recommendation, however, would depend on whether imposing such rules would also prevent borrowers from matching on dimensions that may be desirable from the lender's perspective, such as social connections. If borrowers match with those that they know best, then project covariation is naturally tied to social connectedness because friends or relatives will often have the same income sources and therefore be exposed to similar income shocks. Taken together, endogenous matching will result in groups with both correlated returns and social ties.

In terms of optimal market design, there are three cases to distinguish. First, if the project correlation measure captures social connectedness fully, then group diversification can be implemented by restricting the grouping together of relatives, as suggested in the Grameen Replication Guidelines ([Alam and Getubig, 2010](#)). The remaining two cases are relevant when social connectedness is (partly) captured in the error term. The implications of the findings in this section then also depend on the expected repayment effect of social connectedness. If it improves repayment, pushing for diversification may have no effect (based on the second Probit model in [Table 3](#)). If, on the other hand, social connections have a negative effect on repayment, then there is a clear case for diversifying groups. In the theoretical and empirical literature, there is no clear consensus on the effect of social connections on repayment. For the Thai village context used in this paper, [Ahlin and Townsend \(2007\)](#) find that cooperative behaviour in groups has a negative effect on repayment. This is consistent with the models of [Banerjee](#)

et al. (1994) and Besley and Coate (1995), who predict that cooperation prevents a group from exerting repayment pressure on its members. The result from the survey that most closely matches the context of this paper thus suggests a positive repayment effect from diversification.

6 Conclusion

I analyse the optimal design of rules for group formation in matching markets with an application to group lending in microfinance. The particular focus is on microlenders' decisions on rules to diversify borrower groups with respect to their exposure to common income shocks. Such rules affect group outcomes by influencing who matches with whom (direct effect) and who participates in the market (participation effect). A distinction between these effects allows a direct test of ex-ante and ex-post mechanisms through which the variable of interest affects group outcomes. This distinction is particularly useful in the field of (micro)finance, where the evaluation of adverse selection models requires that moral hazard effects are not in force, and vice versa.

I develop the trade-off for conflicting predictions of extant asymmetric information models and estimate both effects separately. The empirical analysis is complicated by an endogeneity problem that occurs whenever agents match on both (i) the independent variable of interest and (ii) characteristics unobserved to the researcher but correlated with the outcome of interest. To correct for the resulting sorting bias, I develop a generalised Heckman selection model with credible exclusion restrictions that exploits agents' local information to control for unobserved group characteristics. These unobservables are inferred in a matching model that captures the strategic interactions of agents who can only choose from the set of partners that would be willing to match with them.

This paper has implications for empirical and theoretical work on matching markets as well as for microfinance practice, and three main outcomes can be identified. First, empirical studies on group outcomes can correct for bias that results from sorting using R package `matchingMarkets` (Klein, 2015b). Alternatively, empirical findings should be interpreted with this bias in mind, noting that direction and size are often unclear. In the Thai group-lending context in this paper, the positive sorting bias even exceeds the negative direct effect of borrowers' correlated returns on repayment, which has led previous research – using the same dataset – to make incorrect policy recommendations. Second, most theoretical work on microfinance builds on the result that endogenous group formation

is socially optimal when matching is on risk type. Future modelling should take into account that matching also takes place on other dimensions – such as exposure to common shocks – with adverse effects on group repayment. Third, for microfinance practice, this finding suggests that lenders would benefit from ensuring that borrowing groups are sufficiently diversified in their exposure to income shocks. This may be achieved by placing suitable restrictions on the composition of borrower groups.

References

- ABBINK, K., IRLBUSCH, B. and RENNER, E. (2006). Group size and social ties in microfinance institutions. *Economic Inquiry*, **44** (4), 614–628.
- AHLIN, C. (2009). *Matching for credit: Risk and diversification in Thai microcredit groups*. Working Paper 251, Bureau for Research and Economic Analysis of Development.
- and TOWNSEND, R. (2007). Using repayment data to test across models of joint liability lending. *The Economic Journal*, **117** (517), F11–F51.
- ALAM, N. and GETUBIG, M. (2010). *Guidelines for establishing and operating Grameen-style microcredit programs. Based on the practices of Grameen Bank and the experiences of Grameen Trust and Grameen Foundation Partners*. Technical report, Grameen Foundation.
- ALBERT, J. H. and CHIB, S. (1993). Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, **88** (422), 669–679.
- BANERJEE, A., BESLEY, T. and GUINNANE, T. (1994). Thy neighbor's keeper: The design of a credit cooperative with theory and a test. *Quarterly Journal of Economics*, **109** (2), 491–515.
- , DUFLO, E., GHATAK, M. and LAFORTUNE, J. (2013). Marry for what? Caste and mate selection in modern India. *American Economic Journal: Microeconomics*, **5** (2), 33–72.
- BESLEY, T. and COATE, S. (1995). Group lending, repayment incentives and social collateral. *Journal of Development Economics*, **46** (1), 1–18.

- BRESNAHAN, T. and REISS, P. (1991). Empirical models of discrete games. *Journal of Econometrics*, **48** (1-2), 57–81.
- BROCK, W. and DURLAUF, S. (2007). Identification of binary choice models with social interactions. *Journal of Econometrics*, **140** (1), 52–75.
- CASSAR, A. and WYDICK, B. (2010). Does social capital matter? Evidence from a five-country group lending experiment. *Oxford Economic Papers*, **62** (4), 715–739.
- CHE, Y.-K. and YOO, S.-W. (2001). Optimal incentives for teams. *The American Economic Review*, **91** (3), 525–541.
- CHEN, J. (2013). Estimation of the loan spread equation with endogenous bank-firm matching. *Advances in Econometrics*, **3**, 251–289.
- CILIBERTO, F. and TAMER, E. (2009). Market structure and multiple equilibria in airline markets. *Econometrica*, **77** (6), 1791–1828.
- EECKHOUT, J. and MUNSHI, K. (2010). Matching in informal financial institutions. *Journal of the European Economic Association*, **8** (5), 947–988.
- FAFCHAMPS, M. and GUBERT, F. (2007). The formation of risk sharing networks. *Journal of Development Economics*, **83** (2), 326–350.
- GANGOPADHYAY, S., GHATAK, M. and LENSINK, R. (2005). Joint liability lending and the peer selection effect. *The Economic Journal*, **115** (506), 1005–1015.
- GHATAK, M. (1999). Group lending, local information and peer selection. *Journal of Development Economics*, **60** (1), 27–50.
- (2000). Screening by the company you keep: Joint liability lending and the peer selection effect. *The Economic Journal*, **110** (465), 601–631.
- GINÉ, X., JAKIELA, P., KARLAN, D. and MORDUCH, J. (2010). Microfinance games. *American Economic Journal: Applied Economics*, **2** (3), 60–95.
- GORDON, N. and KNIGHT, B. (2009). A spatial merger estimator with an application to school district consolidation. *Journal of Public Economics*, **93** (5-6), 752–765.
- HECKMAN, J. (1979). Sample selection bias as a specification error. *Econometrica*, **47** (1), 153–161.

- HERMES, N. and LENSINK, R. (2007). The empirics of microfinance: What do we know? *The Economic Journal*, **117** (517), F1–F10.
- HOLMSTRÖM, B. and TIROLE, J. (1997). Financial intermediation, loanable funds, and the real sector. *Quarterly Journal of Economics*, **112** (3), 663–691.
- KARLAN, D. (2007). Social connections and group banking. *The Economic Journal*, **117** (517), F52–F84.
- KARLSON, K. B., HOLM, A. and BREEN, R. (2012). Comparing regression coefficients between same-sample nested models using logit and probit: A new method. *Sociological Methodology*, **42** (1), 286–313.
- KATZUR, T. and LENSINK, R. (2012). Group lending with correlated project outcomes. *Economics Letters*, **117** (2), 445–447.
- KLEIN, T. (2015a). *Analysis of stable matchings in R: Package matchingMarkets*. Vignette to R package matchingMarkets, The Comprehensive R Archive Network.
- (2015b). *matchingMarkets: Structural estimators and algorithms for the analysis of stable matchings*. R package version 0.1-5, The Comprehensive R Archive Network.
- KLONNER, S. (2006). *Risky Loans and the Emergence of Rotating Savings and Credit Associations*. Working paper, Cornell University.
- LAFFONT, J.-J. (2003). Collusion and group lending with adverse selection. *Journal of Development Economics*, **70** (2), 329–348.
- MOSLEY, P. (1986). Risk, insurance and small farm credit in developing countries: A policy proposal. *Public Administration and Development*, **6** (3), 309–319.
- PARK, M. (2013). Understanding merger incentives and outcomes in the mutual fund industry. *Journal of Banking and Finance*, **37** (11), 4368–4380.
- PYCIA, M. (2012). Stability and preference alignment in matching and coalition formation. *Econometrica*, **80** (1), 323–362.
- ROTH, A. E. (2008). What have we learned from market design? *The Economic Journal*, **118** (527), 285–310.

- SHARMA, M. and ZELLER, M. (1997). Repayment performance in group-based credit programs in Bangladesh: An empirical analysis. *World Development*, **25** (10), 1731–1742.
- SØRENSEN, M. (2007). How smart is smart money? A two-sided matching model of venture capital. *The Journal of Finance*, **62** (6), 2725–2762.
- STIGLITZ, J. (1990). Peer monitoring and credit markets. *The World Bank Economic Review*, **4** (3), 351–366.
- and WEISS, A. (1981). Credit rationing in markets with imperfect information. *American Economic Review*, **71** (3), 393–410.
- TOWNSEND, R. (1997a). *Townsend Thai Project Initial Bank for Agriculture and Agricultural Cooperatives (BAAC) Survey, 1997*. Tech. rep., Murray Research Archive.
- (1997b). *Townsend Thai Project Initial Household Survey, 1997*. Tech. rep., Murray Research Archive.
- (2000a). *Townsend Thai Project Bank for Agriculture and Agricultural Cooperatives (BAAC) Annual Resurvey, 2000*. Tech. rep., Murray Research Archive.
- (2000b). *Townsend Thai Project Household Annual Resurvey, 2000*. Tech. rep., Murray Research Archive.
- WEESE, E. (2015). Political mergers as coalition formation: An analysis of the Heisei municipal amalgamations. *Quantitative Economics*, **6** (2), 257–307.
- WYDICK, B. (1999). Can social cohesion be harnessed to repair market failures? Evidence from group lending in Guatemala. *The Economic Journal*, **109** (457), 463–475.
- ZELLER, M. (1998). Determinants of repayment performance in credit groups: The role of program design, intragroup risk pooling, and social cohesion. *Economic Development and Cultural Change*, **46** (3), 599–620.
- ZELLNER, A. (1986). *On assessing prior distributions and Bayesian regression analysis with g-prior distributions*, North-Holland, Amsterdam, vol. 6, pp. 233–243.

A Proofs

Proof of Proposition 3.1. Denote by \tilde{p} the *average* success probability of borrowers with risk type $p \in [\underline{p}, \hat{p}]$ who would take a loan at contract terms (r, q) and form groups with project covariation ϵ .

$$\tilde{p} = \frac{\int_{\underline{p}}^{\hat{p}} s g(s) ds}{G(\hat{p})} \quad (\text{A1})$$

This is the expression for the expectation of a truncated distribution with probability density function $g(\cdot)$ and cumulative distribution function $G(\cdot)$. Making use of the selection equation Eqn 2, the expected repayment \tilde{y} of this borrower pool can be written as

$$\tilde{y} = r \frac{\int_{\underline{p}}^{\hat{p}} s g(s) ds}{G(\hat{p})} + q \frac{\int_{\underline{p}}^{\hat{p}} s(1-s) g(s) ds}{G(\hat{p})} - q\epsilon \quad (\text{A2})$$

$$= (r+q) \frac{\int_{\underline{p}}^{\hat{p}} s g(s) ds}{G(\hat{p})} - q \frac{\int_{\underline{p}}^{\hat{p}} s^2 g(s) ds}{G(\hat{p})} - q\epsilon. \quad (\text{A3})$$

Using Leibniz integral rule, quotient rule and the fact that $\int_{\underline{p}}^{\hat{p}} s^2 g(s) ds = (\tilde{p}^2 + \tilde{\sigma}_p^2)G(\hat{p})$, where $\tilde{\sigma}_p^2$ is the variance of the success probability in the borrower pool, we can write the marginal effect of project covariation on expected repayment as

$$\frac{\partial \tilde{y}}{\partial \epsilon} = (r+q) \frac{\hat{p}g(\hat{p})}{G(\hat{p})} \left(1 - \frac{\tilde{p}}{\hat{p}}\right) \frac{\partial \hat{p}}{\partial \epsilon} - q \frac{\hat{p}^2 g(\hat{p})}{G(\hat{p})} \left(1 - \frac{\tilde{p}^2 + \tilde{\sigma}_p^2}{\hat{p}^2}\right) \frac{\partial \hat{p}}{\partial \epsilon} - q. \quad (\text{A4})$$

From Eqn 2 we know that $\partial \hat{p} / \partial \epsilon = q / [r + q(1 - 2\hat{p})]$. Substituting, setting $\underline{p} = 0$ (without loss of generality) and assuming p to be from a uniform distribution⁷ yields

$$\frac{\partial \tilde{y}}{\partial \epsilon} = \frac{1}{2}(r+q) \frac{q}{r+q(1-2\hat{p})} - \frac{2}{3}q\hat{p} \frac{q}{r+q(1-2\hat{p})} - q \quad (\text{A5})$$

$$= \frac{1}{6}q \left[\frac{2q\hat{p}}{r+q(1-2\hat{p})} - 3 \right] < 0 \Leftrightarrow \hat{p} < \frac{3}{8} \frac{q+r}{q}. \quad (\text{A6})$$

This implies that project covariation strictly *reduces* expected repayment if either (i) $\hat{p} < 3/4$ or (ii) $q/r < 3/5$. Consider these results one at a time. For (i), note that, for $q > 0$, $\partial \tilde{y} / \partial \epsilon$ is strictly increasing in joint liability payment q which is

⁷This implies that $\tilde{p} = \frac{1}{2}(\hat{p} + \underline{p}) = \frac{1}{2}\hat{p}$, $\tilde{\sigma}_p^2 = \frac{1}{12}(\hat{p}^2 - \underline{p}^2) = \frac{1}{12}\hat{p}^2$, $g(\hat{p}) = 1/[1 - \underline{p}] = 1$, and $G(\hat{p}) = \frac{\hat{p} - \underline{p}}{1 - \underline{p}} = \hat{p}$.

bounded from above at r . It therefore suffices to analyse the case where $q = r$ for which straightforward calculation (using Eqn A6) results in $\partial\tilde{y}/\partial\epsilon < 0 \Leftrightarrow \hat{p} < 3/4$. Similarly, for (ii), since $\partial\tilde{y}/\partial\epsilon$ is increasing in \hat{p} it suffices to state the condition for \hat{p} close to 1.⁸ In this case, we have $\partial\tilde{y}/\partial\epsilon < 0 \Leftrightarrow q/r < 3/5$. \square

Proof of Corollary 3.1. The proof of this corollary follows directly from Eqn A4 in the proof of Proposition 3.1. The cross partial derivative $\frac{\partial}{\partial g(\hat{p})} \left(\frac{\partial\tilde{y}}{\partial\epsilon} \right) = \frac{\partial^2\tilde{y}}{\partial g(\hat{p})\partial\epsilon}$ is positive if

$$(r + q) \frac{\hat{p}}{\frac{\partial G(\hat{p})}{\partial g(\hat{p})}} \left(1 - \frac{\tilde{p}}{\hat{p}} \right) \frac{\partial\hat{p}}{\partial\epsilon} > q \frac{\hat{p}^2}{\frac{\partial G(\hat{p})}{\partial g(\hat{p})}} \left(1 - \frac{\tilde{p}^2 + \tilde{\sigma}_p^2}{\hat{p}^2} \right) \frac{\partial\hat{p}}{\partial\epsilon} \quad (\text{A7})$$

$$(r + q) \left(1 - \frac{\tilde{p}}{\hat{p}} \right) > q\hat{p} \left(1 - \frac{\tilde{p}^2 + \tilde{\sigma}_p^2}{\hat{p}^2} \right) \quad (\text{A8})$$

$$(r + q)(\hat{p} - \tilde{p}) > q(\hat{p}^2 - \tilde{p}^2) - q\tilde{\sigma}_p^2. \quad (\text{A9})$$

It can be checked that, for $q \leq r$ and $\hat{p} > \tilde{p}$, it holds that $(r + q)(\hat{p} - \tilde{p}) > q(\hat{p}^2 - \tilde{p}^2)$ and therefore the above inequality is satisfied for all parameter constellation in Ghatak (1999). The condition $q \leq r$ is an incentive compatibility constraint. The rationale behind this constraint is that if joint-liability q were to exceed interest payment r , the borrower with the successful project would prefer to announce success and pay interest $r < q$ instead of the full joint-liability payment (Gan-gopadhyay *et al.*, 2005). \square

Proof of Proposition 3.2. The starting point of the proof are two identical, hazardous projects L and M between which borrowers are indifferent.

$$V_{L-M} = V_L - V_M \quad (\text{A10})$$

$$= [p_H^2 + \epsilon] \cdot U_H + [p_H(1 - p_H) - \epsilon] \cdot U_{Hq} \quad (\text{A11})$$

$$-[p_H^2 + \epsilon] \cdot U_H - [p_H(1 - p_H) - \epsilon] \cdot U_{Hq} = 0.$$

Now consider an increase in ϵ for both projects. How much safer can the first project be made in response when (i) the risk-return ratio is fixed at dy/dp and (ii) the borrowers are to be held indifferent between the safer and the risky project? Taking the total differential with respect to ϵ for both projects and allowing a

⁸Note that for $\hat{p} = 1$ we have $\partial\hat{p}/\partial\epsilon = 0$ because $p \in [p, 1]$ and thus $\partial\tilde{y}/\partial\epsilon = -q < 0$ from Eqn A4.

simultaneous change in p and y for the first project yields:

$$\begin{aligned}
dV_{L-M} &= (U_H - U_{Hq}) \cdot d\epsilon + [(U_H - U_{Hq}) \cdot 2p_H + U_{Hq}] \cdot dp \\
&\quad + [(p_H^2 + \epsilon) \cdot U'_H + (p_H(1 - p_H) - \epsilon) \cdot U'_{Hq}] \cdot \frac{dy}{dp} \cdot dp \\
&\quad + (U'_H - U'_{Hq}) \cdot dy \cdot d\epsilon + [(U'_H - U'_{Hq}) \cdot 2p_H + U'_{Hq}] \cdot dy \cdot dp \\
&\quad - (U_H - U_{Hq}) \cdot d\epsilon, \tag{A12}
\end{aligned}$$

where $U'_k = \partial U_k / \partial y$. Setting $dV_{L-M} = 0$ holds the borrower indifferent between the two projects and yields the rate by which an increase in correlation results in a safer project choice, for given level of dy and risk-return ratio dy/dp .

$$\begin{aligned}
dp/d\epsilon &= \left\{ - (U'_H - U'_{Hq}) dy \right\} / \left\{ (U_H - U_{Hq}) 2p_H + U_{Hq} + [(U'_H - U'_{Hq}) 2p_H + U'_{Hq}] dy \right. \\
&\quad \left. + [(U'_H - U'_{Hq})(p_H^2 + \epsilon) + U'_{Hq} p_H] \frac{dy}{dp} \right\}.
\end{aligned}$$

The expected repayment to the bank is

$$Y = r \cdot p_H + q \cdot [p_H(1 - p_H) - \epsilon]. \tag{A13}$$

Taking the total differential w.r.t. p and ϵ yields

$$dY = (r + q(1 - 2p_H)) \cdot dp - q \cdot d\epsilon \tag{A14}$$

$$\frac{dY}{d\epsilon} = (r + q(1 - 2p_H)) \cdot \frac{dp}{d\epsilon} - q. \tag{A15}$$

Substituting $dp/d\epsilon$ from Eqn A13 above into Eqn A15 gives the marginal repayment effect of correlated returns as

$$\begin{aligned}
dY/d\epsilon &= \left\{ (r + q(1 - 2p_H))(U'_{Hq} - U'_H) dy \right\} / \left\{ (U_H - U_{Hq}) 2p_H + U_{Hq} \right. \\
&\quad \left. + [(U'_H - U'_{Hq}) 2p_H + U'_{Hq}] dy + [(U'_H - U'_{Hq})(p_H^2 + \epsilon) + U'_{Hq} p_H] \frac{dy}{dp} \right\} - q.
\end{aligned}$$

Observe that there are two situations in which the marginal repayment effect is strictly negative. First, if borrowers are risk neutral or moderately risk averse such that $U'_{Hq} \approx U'_H$ then $dY/d\epsilon = -q < 0$. In this case, utility is close to linear and correlation has no effect on decision between projects but a strictly negative effect from anti-diversification. The second case is when dy goes towards zero. Then $dY/d\epsilon = -q < 0$ because the income level at which the utility gain from avoiding liability payment (due to increased project correlation) is evaluated – and thus

the slope of the utility – is similar for safe and risky projects. \square

Proof of Proposition 3.3. Part (i) of the proposition is trivial. For part (ii), note that the matching pattern in Figure 2c is an equilibrium under aligned preferences if the safest risk type in group L , denoted by k , prefers to remain matched with group member j' over a swap of j' for borrower i' from group M , i.e. if $u_{k,j'} > u_{k,i'}$. This is the case for $\epsilon > p_k(p_{i'} - p_{j'})$. If this inequality holds for marginal type k , then any borrower x in L prefers j' over i' (because $p_x < p_k$ and net utility in Eqn 5 is decreasing in p). Thus, preferences are aligned within the leading exposure type A for the dominant group L and the matching is stable.

The condition for aligned preferences is satisfied if exposure intensity ϵ is sufficiently large and the difference $p_{i'} - p_{j'}$ is sufficiently small. The latter term is decreasing in the proportion of the leading exposure type. To see this, note that the integral over the probability density function of risk type p below must equal 1/2 in the case with two groups per market.

$$\theta_A \int_{p_{i'}}^{p_{j'}} f_p(t) dt = \frac{1}{2} \quad (\text{A16})$$

Here, the integral is pre-multiplied with the proportion of A -types, θ_A , because, by Assumption H1, the proportion is constant for any point in the distribution of p . Now, fix any distribution of risk types, f_p , and note that the higher the proportion of A -types, θ_A , in the market the higher the value of $p_{j'}$, the lowest risk type in group L . Thus the smaller is the difference $p_{i'} - p_{j'}$. Graphically, in Figure 2b, for any distribution of risk types, the more A -type borrowers, the smaller the term $p_{i'} - p_{j'}$. \square

Proof of Corollary 3.2. To begin with, under Assumption H1 both groups M and L have the same group project correlation and L has safer risk types than M (see Figure 2b). An i -for- j swap has two effects.

First, it results in an increase in project covariation for group L and a decrease for group M . To see this, note that the total differential of Eqn 6 with respect to n_s is $q\epsilon \sum_{s \in \{A,B\}} (2n_s - 1) dn_s$. For an i -for- j swap in group L we have $dn_A = +1$ and $dn_B = -1$, which results in an *increase* in group project correlation of $2q\epsilon(n_A - n_B) > 0$, where the sign of the inequality results from the fact that $n_A > n_B$. Conversely, for group M , setting $dn_A = -1$ and $dn_B = +1$ we observe a *decrease* in group project correlation by $2q\epsilon(n_B - n_A) < 0$.⁹

⁹After several A -for- B swaps, group M may eventually have more B -types than A -types, i.e. $n_B^M > n_A^M$ (see Figure 2c) and project correlation increases. However, the correlation of L still

Second, such a swap increases the average riskiness of types in group L but never makes L riskier than M on average. It follows that sorting induces a positive correlation between the two dimensions. \square

Proof of Proposition 3.4. A matching is stable if deviation is unattractive. Alternative matches are therefore bound to have a lower valuation than observed ones. Specifically, the valuation of an *unmatched* group G must be smaller than the maximum valuation of the equilibrium matches $\mu(i)$ that its members i belong to. If G 's valuation was larger, then its members would block their equilibrium matches to form the new coalition G . We thus have an *upper bound* \overline{V}_G for the valuation of $G \notin \tilde{\mu}$.

$$G \notin \tilde{\mu} \Leftrightarrow V_G < \max_{i \in G} V_{\mu(i)} =: \overline{V}_G \quad (\text{A17})$$

For the *if* direction (\Rightarrow) assume for contradiction that G is a *blocking coalition* for μ . Per the definition of blocking coalitions, this implies that all agents in this coalition prefer being matched to each other over being matched to their current partners in μ , i.e., $G \succ_i \mu(i) \forall i \in G$. Given aligned preferences, the condition implies that $V_G > V_{\mu(i)} \forall i \in G$. Together this implies that $V_G > \max_{i \in G} V_{\mu(i)}$, which contradicts the assumption in the proposition.

For the *only if* direction (\Leftarrow) assume μ to be a stable matching with $G \notin \mu$. Since by stability G is not a blocking coalition, it must hold that there is at least one individual i that prefers its equilibrium group $\mu(i)$ over group G , i.e. $\exists i \in G : \mu(i) \succ_i G$. Given aligned preferences, this condition implies that $\exists i \in G : V_{\mu(i)} > V_G$. Together these conditions imply that $V_G < \max_{i \in G} V_{\mu(i)}$, which is the upper bound condition from the proposition.

Following the same logic as above, the valuation of a *matched* group G must be larger than the maximum valuation of the feasible deviations of its group members. Feasible deviations of G 's group members are such that they are attractive to those borrowers outside of group G that are necessary to form these new matches. That is, feasible deviations are such that their value is larger than the maximum valuation of the equilibrium groups that the non-group- G members of that deviating group belong to.

$$G \in \tilde{\mu} \Leftrightarrow V_G > \max_{G'' \in S} V_{G''} =: \underline{V}_G \quad (\text{A18})$$

Here, S is the set of feasible deviations from G , defined as $S(G) := \{G' \in$

grows at a faster rate.

$\mathcal{G} \setminus \{G' \cap G \notin \{\emptyset, G\}, V_{G'} > \max_{i \in G' \setminus G} V_{\mu(i)}\}$. That is, a deviation from G to G' is feasible for all new non- G borrowers in G' if the valuation of G' is larger than the maximum that new borrowers would receive in their equilibrium match, i.e. if $V_{G'} > \max_{i \in G' \setminus G} V_{\mu(i)}$. The set of new borrowers are those borrowers in G' that do not belong to the original equilibrium match G , i.e. those in $G' \setminus G$.

For the *only if* direction (\Leftarrow) assume μ to be a stable matching with $G \in \mu$. Since μ is stable, no member of G can benefit from deviating. Given aligned preferences, for any member $i \in G$ this implies that $V_G > V_{G'} \forall G' \in S$, where S is the set of feasible deviations for group members of G . Together this implies the inequality $V_G > \max_{G' \in S} V_{G'}$ in the proposition.

For the *if* direction (\Rightarrow) assume that the inequalities in the proposition hold. Let G be a match in μ . It follows from the inequalities in the proposition that no member of G can be part of a blocking coalition. \square

B Simulation of posterior distribution

The Bayesian estimator uses the data augmentation approach (proposed by [Albert and Chib, 1993](#)) that treats the latent outcome and valuation variables as nuisance parameters.

Conditional posterior distribution of outcome variables

The outcome equation is defined (and observed) for realised matches, $G \in \mu$, only. For binary outcome variables, when the observed outcome Y_G equals one, the conditional distribution of the latent outcome variable Y_G^* is truncated from below at zero as $N(X_G\beta + (V_G - W_G\alpha)\delta, 1)$ with density

$$\begin{aligned} \mathbb{P}(Y_G^*|V, Y_{-G}^*, \theta, Y, \mu, W, X) &= C \cdot \mathbb{1}[Y_G^* \geq 0] \\ &\cdot \exp\left\{-0.5(Y_G^* - X_G\beta - (V_G - W_G\alpha)\delta)^2\right\}. \end{aligned}$$

When Y_G equals zero, the distribution is the same but now truncated from *above* at zero. In markets with one group only, the term $(V_G - W_G\alpha)\delta$ is dropped because V_G , α and δ need not be estimated in this case. When an offset is used in the estimation, the distributions are truncated at minus the group-specific offset value instead of zero.

Conditional posterior distribution of valuation variables

For unobserved matches, $G \notin \mu$, the distribution of the latent valuation variable is $N(W_G\alpha, 1)$, truncated from above at \overline{V}_G with density

$$\begin{aligned} \mathbb{P}(V_G|V_{-G}, Y^*, \theta, Y, \mu, X, W) &= C \cdot \mathbb{1}[V_G \leq \overline{V}_G] \\ &\cdot \exp\left\{-0.5(V_G - W_G\alpha)^2\right\}. \end{aligned}$$

For observed matches, $G \in \mu$, the conditional distribution of the latent valuation variable is truncated from below at \underline{V}_G as $N(W_G\alpha + (Y_G^* - X_G\beta)\delta/(\sigma_\xi^2 + \delta^2), \sigma_\xi^2/(\sigma_\xi^2 + \delta^2))$ with density

$$\begin{aligned} \mathbb{P}(V_G|V_{-G}, Y^*, \theta, Y, \mu, X, W) &= C \cdot \mathbb{1}[V_G \geq \underline{V}_G] \cdot \exp\left\{-0.5\left(V_G \right. \right. \\ &\quad \left. \left. - W_G\alpha - \frac{(Y_G^* - X_G\beta)\delta}{\sigma_\xi^2 + \delta^2}\right)^2 \cdot \frac{\sigma_\xi^2 + \delta^2}{\sigma_\xi^2}\right\}. \end{aligned}$$

The variance of $\sigma_\xi^2/(\sigma_\xi^2 + \delta^2)$ for the valuation variables is chosen such that the variance of the error term in the selection equation, σ_η^2 , equals one.¹⁰

Conditional posterior distribution of parameters

Alpha

The coefficient vector α in the selection equation is only estimated for the subset of markets with two borrower groups. This subset is denoted by T_2 and, together with the set of one-group markets T_1 , makes the total set of markets T . The conditional posterior of α is $N(\hat{\alpha}, \hat{\Sigma}_\alpha)$, where

$$\hat{\Sigma}_\alpha = \left[\Sigma_\alpha^{-1} + \sum_{t \in T_2} \left[\sum_{G \notin \mu_t} W'_G W_G + \sum_{G \in \mu_t} \frac{\sigma_\xi^2 + \delta^2}{\sigma_\xi^2} W'_G W_G \right] \right]^{-1} \quad (\text{A19})$$

and

$$\begin{aligned} \hat{\alpha} = & -\hat{\Sigma}_\alpha \left[-\Sigma_\alpha^{-1} \bar{\alpha} + \sum_{t \in T_2} \left[\sum_{G \notin \mu_t} -W'_G V_G \right. \right. \\ & \left. \left. + \sum_{G \in \mu_t} \frac{\sigma_\xi^2 + \delta^2}{\sigma_\xi^2} W'_G \left(V_G - \frac{(Y_G^* - X_G \beta) \delta}{\sigma_\xi^2 + \delta^2} \right) \right] \right] \end{aligned} \quad (\text{A20})$$

The variables Σ_α^{-1} and $\Sigma_\alpha^{-1} \bar{\alpha}$ are constants given the priors. In the estimation, I chose the priors $\bar{\alpha} = 0_{|\alpha|,1}$ and $\Sigma_\alpha = 10 \cdot I_{|\alpha|}$, where $0_{n_1, n_2}$ is the zero matrix of dimension $n_1 \times n_2$ and I_n is the identity matrix of dimension n . The values of the two constants are therefore $\Sigma_\alpha^{-1} = (10 \cdot I_{|\alpha|})^{-1}$ and $\Sigma_\alpha^{-1} \bar{\alpha} = 0_{|\alpha|, |\alpha|}$ respectively.

Beta

Similarly, the conditional posterior distribution of β is $N(\hat{\beta}, \hat{\Sigma}_\beta)$, where

$$\hat{\Sigma}_\beta = \left[\Sigma_\beta^{-1} + \sum_{t \in T_1} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X'_G X_G + \sum_{t \in T_2} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X'_G X_G \right]^{-1} \quad (\text{A21})$$

¹⁰ $\sigma_\eta^2 = \text{var}\left(\frac{\varepsilon \delta}{\sigma_\xi^2 + \delta^2} + x\right) = \frac{(\sigma_\xi^2 + \delta^2) \delta^2}{(\sigma_\xi^2 + \delta^2)^2} + \sigma_x^2 = \frac{\delta^2}{(\sigma_\xi^2 + \delta^2)} + \sigma_x^2$. So $\sigma_\eta^2 = 1$ iff $\sigma_x^2 = \sigma_\xi^2 / (\sigma_\xi^2 + \delta^2)$.

and

$$\hat{\beta} = -\hat{\Sigma}_\beta \left[-\Sigma_\beta^{-1} \bar{\beta} - \sum_{t \in T_1} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X'_G Y_G^* - \sum_{t \in T_2} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X'_G (Y_G^* - \delta(V_G - W_G \alpha)) \right]. \quad (\text{A22})$$

Here, the values of the two constants are $\Sigma_\beta^{-1} = (10 \cdot I_{|\beta|})^{-1}$ and $\Sigma_\beta^{-1} \bar{\beta} = 0_{|\beta|, |\beta|}$ respectively.

Delta

Finally, for δ the posterior is $N(\hat{\delta}, \hat{\sigma}_\delta^2)$, with

$$\hat{\sigma}_\delta^2 = \left[\frac{1}{\sigma_\delta^2} + \sum_{t \in T_2} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} (V_G - W_G \alpha)^2 \right]^{-1} \quad (\text{A23})$$

and

$$\hat{\delta} = -\hat{\sigma}_\delta^2 \left[-\frac{\bar{\delta}}{\sigma_\delta^2} - \sum_{t \in T_2} \sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} (Y_G^* - X_G \beta)(V_G - W_G \alpha) \right]. \quad (\text{A24})$$

Analogously, the values of the two constants are $\frac{1}{\sigma_\delta^2} = \frac{1}{10}$ and $\frac{\bar{\delta}}{\sigma_\delta^2} = 0$.

C Replication Guide¹¹

The results reported herein are fully replicable using the `knitr` literate programming engine in the R open-source software environment for statistical computing. R packages used are: `foreign`, `knitr`, `matchingMarkets`, `reshape`, `survival`, `tseries`.

C.1 Data sources and preparation

All files for replication are in the `inputs/` folder. Documentation and original data used in the paper are in `inputs/rawdata/` and can be directly downloaded in zip format from the Harvard Dataverse:

- 1997 BAAC survey (study_10676) at <http://hdl.handle.net/1902.1/10676>
- 1997 Household survey (study_10672) at <http://hdl.handle.net/1902.1/10672>
- 2000 BAAC survey (study_12057) at <http://hdl.handle.net/1902.1/12057>
- 2000 Household survey (study_10935) at <http://hdl.handle.net/1902.1/10935>

These files are preprocessed using the script in `code/1-0-data-preparation.R` and the cleaned and transformed data is written to the `inputs/data/` folder for analysis.

C.1.1 Group-level variables

I start the preprocessing with the 1997 group-level data in Ahlin and Townsend (2007). This data is not used in the analysis because it lacks individual-level information. It serves two purposes: First, it allows me to verify the correct implementation of the variable transformations in Ahlin and Townsend (2007) which are subsequently applied to the 2000 group-level data in this paper. Second, information on the borrower group age in the 1997 data is used to impute this missing variable in the 2000 data.

C.1.2 Regression imputation of group age

The imputation proceeds in three steps. In the first step, a regression model that explains group age is estimated. This model combines data from the Bank for Agriculture and Agricultural Cooperatives (BAAC) 1997 and 2000 surveys in an

¹¹This section of the Appendix is not intended for publication.

interval regression. While the group age is not observed in the BAAC 2000 data, the quasi-panel still allows me to find bounds for a group’s age (see Table A1 for a summary). Note first that groups from villages that only had a single group in the BAAC 1997 can be no older than this group’s age in the BAAC 1997 survey plus three. Furthermore, for all other villages we know that the log-age of groups in the BAAC 2000 survey can be no larger than 34 (= 2000 – 1966) because the BAAC started its group lending operations in 1966. Finally the BAAC 2000 contains a group history of events such as the admission of new members or the assistance members provided to their peers. The first event documented in this history sets a lower bound on a group’s age, which is otherwise bounded from below at 1.

Table A1: Definition of bounds for interval regression of the missing group_age variable

Groups from	lower bound	upper bound
<i>BAAC 1997 survey</i>	group_age ₉₇	group_age ₉₇
<i>BAAC 2000 survey</i>		
- in villages with single group in '97	$\max\{\text{group_hist}_{00}, 1\}$	$\max_age_{97}+3$
- in all other villages	$\max\{\text{group_hist}_{00}, 1\}$	2000-1966

The results of the interval regression are presented in Table A2 below. The independent variables are explained in Table 2. PCG_membership is a village-level variable that gives the percentage of the village population that is a member of a production credit group. Intuitively, we would expect to find less mature groups in a village were PCG membership is prevalent because this may indicate that BAAC operations in that village started more recently. The expected effect of other variables follows similar reasoning. For example, both group size and loan size are expected to be associated with higher group age simply because groups tend to attract new members as they mature and the loan size typically increases for more mature borrowers.

In the second step, the model above is used to predict the group age for groups in the BAAC 2000 data. In the final step, the uncertainty is reintroduced into the imputations by adding the prediction error into the regression. This is done by adding the working residuals of the interval regression model to the predicted values. The result is plotted in Figures A1a and A1b below where the predicted values are on the straight line; dots represent the original BAAC 1997 data and circles depict the imputed data.

The validity of the imputations is tested by comparing the imputed `group_age` to the upper and lower bounds in Table A1. The fact that the predictions remain well within the bounds for *all* 68 groups in the BAAC 2000 data gives us some confidence in the model.

C.1.3 Borrower-level variables

Borrower-level variables are constructed based on the 2000 BAAC survey and the combined borrower and group level data is in `data/borrower-level.RData`.

C.1.4 Matching data

The core part of the data preparation is the generation of group characteristics based on borrower-level variables for both factual and counterfactual groups. This is implemented and documented in function `stabit` in R package `matchingMarkets` (Klein, 2015b). The resulting group-level data is in `data/group-level.RData`.

C.2 Descriptive statistics, models and simulations

The R code in `inputs/code/` for descriptive statistics, econometrics and simulation results is commented and can be run independently to obtain all results in figures, tables and text in the paper. The code is annotated with tags of the form `## ---- label:`, which allow the identification of the section in `sections/` that a code chunk is called from in the L^AT_EX document. To see how results from the R code are embedded in the paper, see the `.Rnw` files whose file names correspond directly to the tag of the code chunk in the R script.

The estimator developed in this paper is implemented in the R package and the source code available on the Comprehensive R Archive Network. To test the functionality of the software implementation in this package, Klein (2015a) provides simulation evidence of the correct implementation of both design matrix generation and estimators.

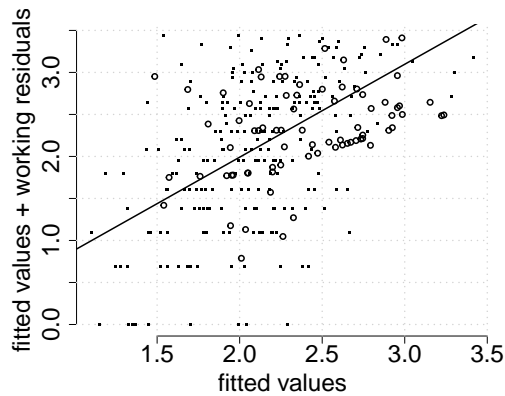
Table A2: Interval regression imputation of the missing group_age variable

*S.E. in parentheses; significance at 0.1, 1, 5, 10% denoted by ***, **, *, and . respectively.*

Interval regression		
<i>Dependent variable as defined in Table A1.</i>		
Intercept	1.451 (0.497)	**
ln(group_size)	0.871 (0.118)	***
loan_size	0.005 (0.006)	.
loan_size_sqrd	-0.000 (0.000)	.
average_land	0.007 (0.002)	**
average_education	-0.548 (0.135)	***
PCG_membership	-0.631 (0.276)	*
BAAC 2000 (ref: 1997)	0.371 (0.125)	**
ln(scale)	-0.332 (0.043)	***
Observations	306	
R^2	0.245	
LR-test, $Pr(> \chi_7^2)$	1e-14	

Figure A1: Comparison of distributions of original group_age variable in BAAC 1997 (dots) and random regression imputation of missing BAAC 2000 group_age variable (circles)

(a) Actual observations (dots) and regression imputations (circles) plotted against fitted values



(b) Actual residuals (dots) and imputed residuals (circles) plotted against fitted values

