# DISCUSSION PAPER

// THILO KLEIN

## Matching for Credit: Identifying Information Asymmetries in Joint-Liability Lending

Leibniz Association

ZEW

# Matching for Credit: Identifying Information Asymmetries in Joint-Liability Lending[*]

Thilo Klein[†]

October 31, 2024

### Abstract

Microcredit, a financial tool providing uncollateralized loans to low-income individuals, has seen a shift from joint-liability (JL) to individual liability (IL) lending models. This article tests a theory explaining this shift, focusing on borrowers matching into groups exposed to similar economic shocks under JL, diminishing its effectiveness. I reconcile conflicting theoretical predictions and propose an empirical strategy to distinguish adverse selection from moral hazard effects. Using data from Thailand, I find that increasing diversity within borrower groups leads to a 10 percentage point improvement in timely repayment. These results inform contract design and strategies to reduce information asymmetries in lending practices.

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

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

**Motivation.** This article investigates successful contractual arrangements used by real-world lending institutions, focussing on group-lending schemes and credit cooperatives in developing countries. These contracts, which delegate screening, monitoring, and enforcement to self-selected group members, have significantly expanded financial access for over 200 million borrowers, particularly low-income households without collateral (Reed, 2015). However, their effectiveness diminishes in rural areas, where financial exclusion is most severe. The high correlation of agricultural loans poses a barrier to lending in these markets (Mosley, 1986).

---

[†]Pforzheim University and ZEW – Leibniz Centre for European Economic Research, Mannheim, Germany. E-mail: thilo.klein@zew.de

Institutions like the Grameen Bank have expressed concerns about joint-liability lending in agriculture, as concurrent project failures among group members lead to collective default. Ahlin (2020) even shows that borrowers choose group members facing similar shocks to avoid joint-liability payments, indicating scope for lender intervention to prevent such groupings.[1]

Understanding the impact of project correlation on group formation, borrowing decisions, project choices, and repayment behaviour is crucial for the optimal credit contract design. This article evaluates the effectiveness of preventing the grouping of borrowers exposed to similar shocks compared to other mechanisms for mitigating information asymmetries. Specifically, it examines the effect of project correlation on adverse selection and moral hazard, providing a basis for comparison with other strategies to mitigate information asymmetries, such as dynamic contracts, enhanced screening strategies, subsidies, loan guarantees, and information coordination. These insights are particularly pertinent given the industry trend toward individual-liability lending (Attanasio *et al.*, 2015) and the increasing extension of loans to non-poor clients in urban areas (de Quidt *et al.*, 2018; Reed, 2015).

**Challenges.** The theoretical literature extensively examines the impact of correlated returns on repayment through direct channels and models incorporating adverse selection and moral hazard. However, a comprehensive analysis integrating these effects to establish conditions under which correlation lowers overall repayment is lacking. For instance, Ghatak (2000) notes that with perfect project correlation, joint-liability contracts perform poorly, as simultaneous project failures lead to bank losses. Conversely, Ahlin and Townsend (2007) show that positive project covariation can enhance repayment by mitigating information asymmetries. This article reconciles the conflicting effects, integrating them into a unified framework.

Empirically testing these theoretical predictions is intricate due to potentially counteracting effects and confounded by a *sample selection bias*, when borrowers sort into groups on both observables and unobservables, as noted in the literature (Hermes and Lensink, 2007). This bias complicates the identification of causal effects and is inherent in coalition formation games in various settings. For example, in a firm's management of intercultural teams, the *(adverse) selection effect* might reduce the worker pool due to preferences against mixed teams, while the

---

[1]Such policies appear feasible. Evidence suggests banks like those operating under the Grameen model explicitly prohibit the grouping of relatives to avoid collusion (see Alam and Getubig, 2010, p. 17).

*treatment effect* considers the impact of team composition on outcomes for a fixed pool of workers. The *sample selection bias* arises if open-minded workers, who naturally perform better, are more likely to join mixed teams. This bias overstates the treatment effect of mixed teams by conflating it with the positive impact of open-mindedness. This article proposes an identification strategy to distinguish treatment and selection effects, empirically testing the theoretical models' predictions.

**This Article.** This article analyses group formation in joint-liability credit groups, focusing on correlations in project returns. The study examines theoretical models of adverse selection and both ex-ante and ex-post moral hazard to identify conditions under which project correlation reduces overall repayment. A structural model is developed and tested using data from agricultural group loans in Thailand. To address the sample selection bias from endogenous group formation, the model incorporates information on potential groups that did not form, similar to the Heckman correction approach. Using counterfactuals, the article evaluates the aggregate welfare effect of rules on permissible group compositions, such as restrictions on grouping individuals exposed to similar economic shocks.

**Model.** This article makes several key assumptions. It primarily uses the Stiglitz and Weiss (1981) model, assuming risk-neutral agents. While this assumption is relaxed in parts of the theoretical analysis, it is maintained throughout the empirical section. The theoretical model further assumes a continuum of borrowers, allowing for sharp predictions applicable to group sizes larger than or equal to two. For the empirical model, this setup is adapted to a finite number of agents. In the empirical analysis, payoffs are divided through Nash bargaining over expected group outcomes, with agents having equal bargaining power. This results in pairwise-aligned preferences (Pycia, 2012), ensuring that if agent $i$ prefers partner $j$ over partner $k$, the joint expected payoff of pair $ij$ is also higher than that of pair $ik$. This preference structure allows summarising the preferences of a coalition into a single match valuation, predicting a unique stable matching, even when outside options are endogenous (Talamas, 2020), thus making the empirical model tractable. Additionally, the error terms in the sample selection model are assumed to follow a bivariate normal distribution, a common assumption in the literature that allows for a clearer exposition of the novel aspects of the estimator derived in this article.

**Data.** The sample selection model is applied to a resurvey of Townsend (2000a) on joint-liability groups of the Bank for Agriculture and Agricultural Co-operatives (BAAC). The BAAC is the largest lender in rural Thailand. The resur-

vey comprises 39 villages from two regions, randomly sampled with stratification. In every village, up to two BAAC groups were surveyed, resulting in data on 68 groups and 316 borrowers. This dataset enables the construction of all possible counterfactual groups to model the endogenous group formation explicitly. Key concepts in this study include borrowers' timely repayment, risk type, exposure to common economic shocks, and the cost of monitoring.

**Method.** The empirical analysis employs a sample selection model to estimate the *treatment effect*, while correcting for *sample selection bias*. This model extends the Heckman (1979) selection correction to a group formation process context, requiring instrumental variables that determine which groups are formed (instrument relevance) but not group outcomes (instrument exogeneity). Identifying suitable instruments is challenging when group sorting is aimed at optimising outcomes. This article demonstrates how the matching interaction itself can provide identification. The equilibrium matching is characterised by bounds on the valuations of observed groups, determined by the outside options of group members to deviate from their group and form new matches. These bounds are a function of the characteristics of other agents in the market. The identification strategy relies on the exclusion restriction that while these characteristics (summarised by the bounds) influence group formation, they do not affect the outcomes of matched groups, which depend solely on the group members. This approach allows for point-identification and estimation of the treatment effect (see Section 1 for an illustrative example).

**Findings.** This article develops the key trade-off of conflicting effects suggested in the literature. The negative repayment effect identified by Ghatak (2000) is dominant for most parameter constellations in both the adverse selection and the ex-ante moral hazard models. Conversely, the ex-post moral hazard model shows a strictly positive effect on repayment, which decreases with higher monitoring costs. Empirical tests of the moral hazard models of Stiglitz (1990) and Armendáriz (1999) support these findings. For ex-ante moral hazard, project covariation significantly reduces timely repayment by 20 percentage points per standard deviation increase in correlation. This negative treatment effect is net of a positive sample selection bias of 27 percentage points, as groups with higher project correlation tend to have better observable and unobservable characteristics. Moderate support is found for the ex-post moral hazard model, with repayment decreasing as monitoring costs rise.

Counterfactual analysis using parameter estimates from the sample selection model tests the aggregate effect of matching on risk exposure. Allowing or pro-

hibiting matching based on exposure type, while keeping other model parameters fixed, reveals that prohibiting the matching on exposure type (diversification) does not draw existing borrowers out of the programme. As predicted, the negative effect of concurrent project failures and loss of joint-liability payment dominates. Preventing matching on exposure type (diversification) improves timely repayment by 10 precentage points.

**Literature.** This article makes two main contributions to the economic literature. First, it advances structural empirical work on matching markets by developing a sample selection model that simultaneously estimates a one-sided matching model and corrects for selection bias in the outcome equation. This extends the work of Sørensen (2007), Chen (2013), and Park (2013), who focus on two-sided markets. Unique to this article is the application of this methodology to one-sided matching market. To ensure the model's likelihood is well-defined (Bresnahan and Reiss, 1991; de Paula, 2013), I build on the Nash bargaining model in Talamas (2020) for pairwise-aligned preferences and (Pycia, 2012) for equilibrium uniqueness, to derive specific equilibrium bounds. The first stage matching model is itself broadly applicable in one-sided matching markets, such as school district mergers and municipal amalgamations. It improves upon previous models by avoiding more restrictive assumptions on agents' preferences, such as pairwise symmetry (Gordon and Knight, 2009) or constraints on permissible coalitions (Weese, 2015).[2]

Second, this article provides empirical evidence on the welfare implications of correlated returns in models of asymmetric information. Previous empirical studies on joint-liability lending, such as Ahlin and Townsend (2007), Wydick (1999), Zeller (1998) and Sharma and Zeller (1997), found mixed effects due to sample selection bias, a challenge recognized in the literature (Hermes and Lensink, 2007). Experimental methods, including those by Karlan (2007), Giné et al. (2010) and Abbink et al. (2006), have been used to test theories of joint-liability lending but fail to account for sorting on specific design-relevant variables like correlated returns.[3] Furthermore, field experiments suffer from attrition if agents cannot be

---

[2] The model proposed in this article also diverges from existing literature on network formation in terms of methodology. Models of network formation, exemplified by Fafchamps and Gubert (2007), do not impose constraints on group size, leading to competition among participants for limited positions. This characteristic complicates empirical analysis, which is a focal point of the present study. Unlike the focused investigations of Klonner (2006) and Eeckhout and Munshi (2010) into two-sided matching scenarios within fixed-sized chit fund groups, this article concentrates on one-sided matching dynamics.

[3] Karlan (2007) makes use of the quasi-random group assignment of microlender FINCA in Peru to estimate the *treatment effect* of social connections. In framed field experiments, Giné et al. (2010) implement a 'partner choice' treatment to estimate the effect of endogenous group

committed to take a loan before knowing their group members.

**Organisation.** This article is organised as follows. Section 1 illustrates the sample selection bias and lays out the identification strategy. Section 2 develops the key trade-off between the conflicting effects suggested in the literature for a continuum economy and characterises the equilibrium matching in an empirical setting with finite markets. Section 3 presents the identification strategy. Section 4 describes the data and presents the results. Section 5 concludes.

# 1    Example of selection bias and identification

This section provides a comprehensive exploration of the bias arising from sorting into groups and illustrates a source of identification. It places a specific emphasis on the specification error that results when variables influencing both group formation and outcomes are unobserved. This is a special case of the measurement error problem that is controlled for in the model in Section 3.

In the context of a credit market, four entrepreneurs, labelled as $b$, $c$, $d$ and $e$, can take loans in groups of two. Each entrepreneur, represented as $i$, has two characteristics: project risk, denoted by the success probability $p_i \in [0,1]$, and project exposure $s_i \in \{A, B\}$ to either of two independent external shocks $A$ and $B$. Notably, exposure is observable, and I represent the vector of exposure types of these four entrepreneurs as $s = (A, A, A, B)$. The vector of project risk $p = (1, 0.5, 0.25, 0.5)$, however, remains unobservable in the available data. Within this credit market, a total of six possible matches emerge, and their repayment is given by the following outcome equation:

$$Y_{ij} \;\; = \;\; \beta_0 + \beta_1 \cdot 1_{s_i = s_j} + \delta \cdot p_i p_j + \xi_{ij}, \tag{1}$$

where, $Y_{ij}$ is the repayment outcome, $1_{s_i = s_j}$ represents whether the group members share the same exposure type, and $p_i p_j$ is the group risk. For simplicity, we maintain $\xi_{ij}$ as a constant of the value zero. Notably, the true parameter values are set at $\beta_0 = 7/24$, $\beta_1 = 0$, and $\delta = 1$. Given that $\beta_1$ is zero, the repayment outcome remains independent of whether the group members share the same exposure type. Therefore, matching with a group member of the same exposure type does not affect the repayment of the group.

---

formation compared to random assignment. Similarly, *selection* and *treatment 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.

We observe the matches $bc$ and $de$, and the outcomes of these matches are depicted by the black points in Figure 1a. The white points represent the outcomes of matches that remain unobserved. Given the unobservable nature of risk type, one can imagine an estimation estimating approach, focusing on a simplified outcome equation:
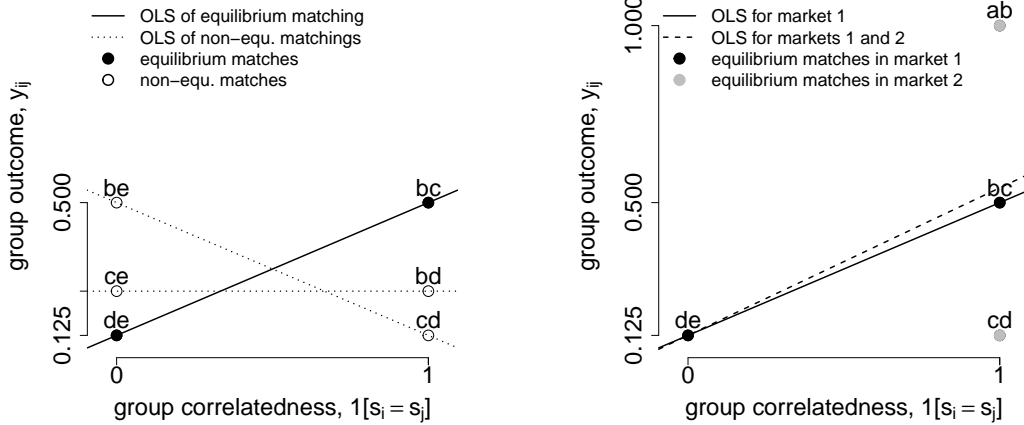
$$Y_{ij} \;=\; \beta_0 + \beta_1 \cdot 1_{s_i = s_j} + \varepsilon_{ij}. \tag{2}$$

In this context, the error term $\varepsilon_{ij}$ captures the unobserved project risk, expressed as $\delta p_i p_j + \xi_{ij}$. The Ordinary Least Squares (OLS) estimate for the slope of Equation 2 is $\hat{\beta}_1 = 0.375$ and it exhibits an upward bias compared to the true value of zero.

Figure 1: Illustration of selection bias and identification in one-sided matching

(a) The regression for equilibrium matching $\mu_1 = \{bc, de\}$ results in an upward-biased slope estimate. The bias resolves (and the slope estimate is zero) when sampling is at random from all three possible matchings.

(b) In a second market, an additional entrepreneur $a$ leads to matching $\mu_2 = \{ab, cd\}$. This facilitates identification, as matches $de$ and $cd$ with similar unobserved characteristics have different correlatedness for exogenous reasons.



The origin of this upwards bias is that $cov(1_{s_i = s_j}, \varepsilon_{ij}) > 0$. That is, match $bc$ has both safer risk type (captured in the error term) and higher correlation than match $de$.

The bias can be effectively resolved through the random assignment of groups. In such a scenario, the slope estimate can be regarded as the equally weighted average of the OLS estimates for three equally probable matchings $\mu_1 = \{bc, de\}$, $\{be, cd\}$ and $\{bd, ce\}$, leading to $\hat{\beta}_1^* = \frac{1}{3} \cdot (0.375 - 0.375 + 0) = 0$.

Furthermore, the bias also be mitigated by introducing an additional market for observation. This new market comprises the same four entrepreneurs, but incorporates a fifth entrepreneur, $a$. In this context, $p_a = 0.95$ and $s_a = A$. The presence of entrepreneur $a$ instigates a shift in the relative ranks within the market, leading to borrower $b$ being matched with $a$, while her former partner $c$ is matched with $d$. Entrepreneur $e$ remains unmatched, resulting in the matching $\mu_2 = \{ab, cd\}$. The outcomes of matchings $\mu_1$ and $\mu_2$ in these two markets are illustrated by the points in Figure 1b. The slope estimate remains upwardly biased, but when comparing groups $de$ and $cd$, it is evident that they have the same repayment outcome, but these groups exhibit differing correlation levels for exogenous reasons. As a result, the outcome remains unaffected by correlation, yielding an estimate of $\beta_1$ as zero.

The presence of borrower $a$ (along with her characteristics) plays a pivotal role in determining who matches with whom and results in the formation of groups with similar outcomes who exhibit differing levels of correlation due to exogenous factors. This allows for meaningful comparison in this example. In the empirical approach outlined in this paper, a matching model is employed to ascertain which groups are comparable across different markets.

## 2   Theoretical framework

This section is divided into two subsections. The first describes the model setup. The second derives the repayment implications of correlated project returns.

### 2.1   Model setup

The model is based on the Stiglitz and Weiss (1981) setting of credit rationing. Project correlation is introduced following Ahlin (2020) and a joint-liability contract in the form of a static liability payment (Ghatak, 1999). 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 $\underline{p} \in (0, 1)$. Agents are of one of three exposure types $s \in \{A, B, N\}$ which constitute the proportions $\theta_A$, $\theta_B$ and $\theta_N$ of the agent population. While $N$-types are not affected by external shocks, the project success probability of $A$- and $B$-types depends on the independent shocks $A$ and $B$, respectively. Specifi-

cally, an $A$-type's probability of success is given by $Pr(y_i > 0) = p_i + \tilde{\gamma}A$, with $A \sim Bin(1, \frac{1}{2})$ coded +1 for 'success' and -1 for 'failure'. The quivalent holds for $B$-types. These shocks equiprobably add probability mass $\epsilon := \tilde{\gamma}^2$ to the symmetric events (where both borrowers succeed or fail) and subtract it from the asymmetric events (where one group member fails and the other succeeds). The expected return $E$ is the same for all risk and exposure types. Under asymmetric information, the lender cannot discriminate between borrower risk types and therefore offers a pooling contract with gross interest rate $r$ and a joint-liability payment $q$ is due in the asymmetric event where borrower $i$ succeeds and partner $j$ fails.

The expected utility of borrower $i$ forming a group with borrower $j$ can then be written as

$$u_{i,j} \quad = \quad E - rp_i - q[p_i(1 - p_j) - \epsilon \cdot 1_{s_i = s_j}]. \tag{3}$$

Here, $1_{s_i = s}$ is an indicator, that is 1 if borrower $i$ is of exposure type $s \in \{A, B\}$ and 0 otherwise and the constant $\epsilon$ gives the intensity of the projects' exposure to shocks. In words, the expected utility is given by the expected project return $E$ less the expected payable interest $rp_i$ and expected joint-liability payment $q[p_i(1 - p_j) - \epsilon \cdot 1_{s_i = s_j}]$. Because agents have no pledgeable collateral, borrower $i$ only pays $q$ in the asymmetric case where her project is successful and partner $j$ defaults.

## 2.2 Revised theories and implications

For the three theoretical models on adverse selection (Ghatak, 1999), ex-ante moral hazrad (Stiglitz, 1990) and ex-post moral hazard (Armendáriz, 1999), 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.

### 2.2.1 Adverse selection

In this setting, Ghatak (1999) shows how the lender can harness joint-liability contracts in groups of two borrowers to mitigate credit rationing. Agents face two decisions: with whom and whether to take a loan. For the first decision, as payoff is super-modular in $p$ and sub-modular in $s$, agents form groups that are homogeneous in both risk and exposure type such that $p_i = p_j$ and $s_i = s_j$. 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

$$E - r\hat{p} - q[\hat{p}(1 - \hat{p}) - \epsilon] \ = \ \bar{u} \tag{4}$$

with equality. Credit is rationed as borrowers with projects safer than $\hat{p}$ do not find it profitable to borrow.

*Developing the key trade-off*

For the effect on repayment, there are two effects to consider. First, Ahlin and Townsend (2007) show 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 $\epsilon$ 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 $p \in [\underline{p}, \hat{p}']$. Second, after an increase in $\epsilon$, the new marginal type $\hat{p}'$ now has the same expected repayment $(E - \bar{u})$ as the previous marginal type $\hat{p}$. However, all inframarginal types now have worse expected repayment (by the term $q \cdot d\epsilon$) because the increase in correlation allows them to avoid liability payments more often. Proposition 2.1 provides conditions for project covariation to reduce repayment when the distribution of risk types is uniform.

**Proposition 2.1.** *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 2.1 is limited to uniform distributions of risk types. Corollary 2.1 below shows that these thresholds are even higher for distributions with lower probability mass in the area of the marginal type.

**Corollary 2.1.** *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 BAAC*

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.[4]

### 2.2.2 Ex-ante moral hazard

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 $k \in \{S, R\}$. Here $S$ is the safe project that was tied to the borrower in the previous subsection and $R$ is the risky project with $p_R < p_S$. The risky project $R$ has a higher *actual* outcome when successful, i.e. $y_R > y_S$, but a lower *expected* outcome, $p_R y_R < p_S y_S$. So the safe project is always socially preferable, but not necessarily privately. 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

---

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

$V_{kk}$, resulting in individual project success probability

$$p \;=\; p_R \cdot 1_{V_{SS} < V_{RR}} \;+\; p_S \cdot 1_{V_{SS} \geq V_{RR}}, \tag{5}$$

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 risky project. The expected group payoff, given project choice $k \in \{S, R\}$, is

$$V_{kk} \;=\; 2 \cdot u(y_k - r) \cdot [p_k^2 + \epsilon] + 2 \cdot u(y_k - r - q) \cdot [p_k(1 - p_k) - \epsilon]. \tag{6}$$

*Developing the key trade-off*

For the effect on repayment, there are again two effects to consider. First, Ahlin and Townsend (2007) show that the utility gain from avoiding joint-liability payment (of size $2q \cdot d\epsilon$) due to an increase in $\epsilon$ is comparatively higher for the safe project, tilting incentives towards choosing the safer project. This is because (i) the safe project has lower returns when successful and (ii) borrowers' utility is concave. Second, there is also a negative effect that correlation has through borrowers avoiding joint-liability payments to the bank. The key trade-off is developed in Proposition 2.2 below.

**Proposition 2.2.** *The marginal effect of project covariation on repayment is strictly negative if either (i) borrowers are risk neutral or (ii) the returns of the risky and safe project are the same.*

*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) $y_R \approx y_S$, the marginal increase in utility from higher project covariation is the same for both the safe and the risky project, $\partial U_{SS}/\partial \epsilon = \partial U_{RR}/\partial \epsilon = 2q$. 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 $\epsilon$ has no effect on project choice. However, it has a strictly negative effect of $-2q \cdot d\epsilon$ from a diversification point of view because it reduces the probability that at least one borrower is successful.

### 2.2.3 Ex-post moral hazard

The ex-post moral hazard model is based on Armendáriz (1999). The difference from the original model is that there is no non-refinancing threat. Such threats are

not common in the BAAC lending programme (Ahlin and Townsend, 2007). The model focuses on whether borrowers monitor each other with the required intensity (1st stage) and whether they subsequently decide to default strategically (2nd stage). The model predicts that positive project covariation maximises the relative benefits from monitoring and lowers the temptation to default strategically. As in the previous models, correlation also has a negative effect from anti-diversification. However, the key trade-off developed in Proposition 2.3 predicts a strictly positive effect for this model.

In the Armendáriz (1999) model, individual project outcome is not common knowledge, but agents can observe their peers' outcome with probability $\gamma$ when investing monitoring effort $\gamma$ at a linear cost $c\gamma$. It is further assumed that the monitoring effort, $\gamma$, is observed by the monitored group member. Strategic default leads to social sanctions by group members which is experienced in the form of a cost $W$. To solve this game, we first find sub-game perfect equilibria in which strategic default does not take place. We then solve the game by backward induction to derive the ex-ante monitoring decision.

**Lemma 2.1.** *Project covariation decreases the minimum monitoring effort necessary to prevent a peer from defaulting strategically.*

*Proof:* See Appendix A.

The intuition of Lemma 2.1 is that positive project covariation increases utility from non-default by putting agents in a position where they are more likely to avoid joint-liability payment, $q$. It therefore makes strategic default less attractive.

**Lemma 2.2.** *Project covariation increases the relative benefits from monitoring a peer.*

*Proof:* See Appendix A.

In summary, positive project covariation maximises the relative benefits from monitoring in the first instance and also lowers the temptation to default strategically in the second instance. In the original paper, these effects are working in opposite directions because it also models a non-refinancing threat.

*Developing the key trade-off*

The following proposition derives the key trade-off between the positive effects derived in Lemmas 2.1 and 2.2 above and the negative effect from borrowers avoiding joint liability payments when their returns are correlated.

**Proposition 2.3.** *The marginal effect of project covariation on repayment is strictly positive.*

*Proof:* See Appendix A.

In contrast to the Ghatak (1999) and Stiglitz (1990) models, the repayment effect of correlated returns is strictly positive in the Armendáriz (1999) model without refinancing threat. Another testable prediction from the model is that this positive repayment effect is strictly decreasing in monitoring cost.

**Corollary 2.2.** *The marginal effect of project covariation on repayment is strictly decreasing in the cost of monitoring.*

*Proof:* See Appendix A.

To see the intuition of Corollary 2.2, recall that, in the model, (i) a borrower can learn about a partner's outcome when investing monitoring effort at a linear cost and (ii) project correlation reduces required monitoring effort, $\gamma$, by making repayment more attractive (see Lemma 2.1). Now, because of the linear cost function $c\gamma$, such cost savings are higher if the increase in correlation is accompanied by a reduction in the unit cost of monitoring. This interaction is tested in the empirical section to identify the ex-post moral hazard effect separately from ex-ante moral hazard.

## 2.3   Characterisation of matchings in finite markets

The theoretical model developed in the previous section provides precise predictions for a continuum of borrowers. Due to the symmetry in equilibrium groups, these predictions also apply to groups with more than two members. For the empirical model, however, this framework must be adapted to accommodate a finite number of agents. In this subsection, I present the modified setup and provide an equilibrium characterisation for the resulting matchings.

### 2.3.1   Model setup

The empirical model considers a finite number of agents forming groups of size $n \geq 2$. In this setup, the utility for borrower $i$ in group $G$ extends the model in Eqn 3 as follows:

$$u_{i,G} \;=\; E - rp_i - qp_i \sum_{j \in G\backslash i}(1 - p_j) + q\epsilon \sum_{s \in \{A,B\}} 1_{s_i = s} \cdot (n_s^G - 1), \qquad (7)$$

where $n_s^G$ represents the count of group $G$ members of exposure type $s$.

Matching in this model involves vertical (risk type) and horizontal (exposure type) dimensions, with super-modular payoffs in risk types (as complements) and sub-modular payoffs in exposure types (as substitutes). In finite markets, stable matching are not guaranteed to exist, as they depend on how utility is shared among agents. Here, agents negotiate binding contracts on the share of future realised payoffs using Nash bargaining (see Altınok, 2023, for a two-player application in joint-liability lending). These contracts are contingent on different states of the world (see Appendix C). Assuming equal bargaining power, the disagreement point is set by the agricultural labour wage, while the outside option for each agent is endogenous, reflecting potential utility from joining alternative groups (Talamas, 2020).[5] Nash bargaining leads to pairwise aligned preferences (Pycia, 2012), such that if agent $i$ prefers partner $j$ to $k$, then joint expected payoff of pair $ij$ is also higher than that of pair $ik$, i.e. $j \succ_i k \Leftrightarrow V_{ij} > V_{ik}$ (Ferdowsian *et al.*, 2020). This preference constraint, widely applied in the literature (Sørensen, 2007; Agarwal, 2015; Dur *et al.*, 2022), enables the aggregation of coalition preferences into a single match valuation, ensuring a unique stable matching, which is essential for the empirical model's tractability.[6]

### 2.3.2 Equilibrium bounds

The equilibrium conditions are represented as inequalities that set lower and upper bounds on the match valuations for both observed and counterfactual (unobserved) matches. In Section 3, I apply these bounds to ensure a unique equilibrium in the empirical model. Proposition 2.4 summarises stability conditions based on bounds $\overline{V_G}$ and $\underline{V_G}$, derived in Appendix A. These conditions hold under aligned preferences, accommodating groups and markets of any size. Although equilibrium groups $G \in \mu$ and non-equilibrium groups $G \notin \mu$ share the same stability The conditions, they impose distinct bounds on the latent valuation variables to identify a unique market equilibrium.

**Proposition 2.4.** *The matching $\mu$ is stable iff $V_G < \overline{V_G} \quad \forall G \notin \mu$. Equivalently, the matching $\mu$ is stable iff $V_G > \underline{V_G} \quad \forall G \in \mu$.*

*Proof:* See Appendix A.

---

[5]While the disagreement point is the agricultural labour wage, outside options are endogenously determined and, per the "outside option principle" (Sutton, 1986), serve as lower bounds on payoffs in bargaining.

[6]A well-defined likelihood for the empirical model requires the observed equilibrium to be a unique stable matching (Bresnahan and Reiss, 1991).

The upper bound $\overline{V_G}$ represents the maximum opportunity cost for each member of group $G$ in leaving their equilibrium group for a non-equilibrium match, while the lower bound $\underline{V_G}$ captures the opportunity cost for members of $G$ to maintain their equilibrium match.
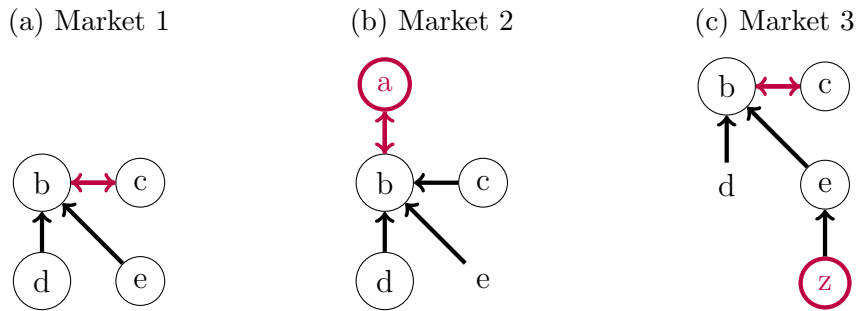
*Example: Equilibrium bounds*

Consider six agents with characteristics outlined in Table 1, who sort into groups of two in three different markets, based on the utility function from Eqn 7 . Market 1, discussed in Section 1, includes agents $b$, $c$, $d$, and $e$.

Table 1: Agent characteristics

| agent | a | b | c | d | e | z |
|---|---|---|---|---|---|---|
| exposure type, $s$ | A | A | A | A | B | B |
| risk type, $p$ | 0.95 | 1 | 0.5 | 0.25 | 0.5 | 0.25 |

In Figure 2a, each agent points to their preferred partner, resulting in at least one stable coalition that is mutually pointed to by its members (Talamas, 2020).[7] In Market 1, this group is $bc$, resulting in the matching $\mu_1 = \{bc, de\}$.

Figure 2: Illustration of matching



(a) Market 1                (b) Market 2                (c) Market 3

The matching $\mu_1$ can be characterised by equilibrium bounds. For unmatched pairs $ij \notin \mu_1$, **upper bounds** are set by the highest opportunity cost of deviating from $\mu_1$ to forming a new match. Here, one of the two observed groups, $bc$ or $de$, is the top coalition to which all group members point (though we do not observe which). Suppose it is $bc$; then, due to pairwise alignment, this group's valuation

---

[7]The unique stable match is identified through a top-down sorting algorithm (Talamas, 2020), where all agents point to their preferred group. One group will always have unanimous member preference; set it aside and repeat until no agents remain.

must exceed that of any counterfactual group containing either $b$ or $c$, such as $cd$. Therefore, the upper bound for the unmatched pair $cd$ is $\overline{V_{cd}} = max\{V_{bc}, V_{de}\}$.

For matched pairs, **lower bounds** are determined by the opportunity cost of maintaining the current match. That is, the valuation of any observed group, say $bc$, must exceed any match where agents from ($bc$) are paired with others who prefer them (to their matched partner in $\mu_1$). Here, $d$ and $e$ both prefer $b$ (see Figure 2a) and would thus like to deviate. They would also prefer to match with $c$. Thus, the lower bound for $bc$ is $\underline{V_{bc}} = max\{V_{bd}, V_{be}, V_{cd}, V_{ce}\}$. For $de$, the lower bound is $\underline{V_{de}} = -\infty$, as neither $b$ nor $c$ prefer deviating from $bc$ to match with $d$ or $e$.

# 3 Empirical strategy

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

## 3.1 Treatment effect

This subsection presents a structural empirical model to estimate the treatment effect of project correlation on repayment while addressing sample selection bias. The observed equilibrium groups constitute a self-selected sample, and the selection problem here differs substantially from the classical two-stage correction of Heckman (1979). Specifically, the first-stage selection mechanism that determines which borrower groups are observed is modeled as a one-sided matching game rather than a simple discrete choice. In discrete choice models, an observed match directly reveals the preferences of group partners toward each other. However, in our context, an observed matching results from complex interactions among agents. Borrowers can only choose partners who are also willing to form a match, but we do not observe their relevant choice sets. This limitation renders direct inference based on discrete choice models infeasible, even when accounting for social interactions as in Brock and Durlauf (2007) and Ciliberto and Tamer (2009).

The empirical strategy, therefore, simultaneously estimates the repayment outcome equation alongside the matching game, expressed by the valuation equation:

$$V_G \;=\; W_G\alpha + \eta_G, \tag{8}$$

where $V \in \mathbb{R}^{|\Omega|}$ is the vector of latent valuations for all feasible groups $\Omega$ in the market, $W \in \mathbb{R}^{|\Omega| \times k}$ is the matrix of $k$ group characteristics, $\alpha \in \mathbb{R}^k$ is a parameter vector, and $\eta \in \mathbb{R}^{|\Omega|}$ is a vector of random errors. There are $|\Omega|$ such equations, corresponding to all feasible groups.[8]

A group – and this its repayment outcome $Y_G$ – is observed if it is part of the equilibrium matching $\mu$, meaning its valuation $V_G$ lies within the set $\Gamma_\mu$ satisfying the equilibrium conditions.[9] This set links the structural empirical model to the theoretical equilibrium characterizations derived in Proposition 2.4, Subsection 2.3. The equilibrium condition can be expressed as:

$$V \in \Gamma_\mu \Leftrightarrow \left[ V_G < \overline{V_G} \ \forall G \notin \mu \right] \Leftrightarrow \left[ V_G > \underline{V_G} \ \forall G \in \mu \right], \tag{9}$$

where $\overline{V_G}$ and $\underline{V_G}$ are the upper and lower bounds on valuations for unmatched and matched groups, respectively.

The outcome equation for the binary dependent variable is defined as $Y_G = 1[Y_G^* > 0]$, with the latent outcome $Y_G^*$ given by:

$$Y_G^* \;=\; X_G \beta + \varepsilon_G. \tag{10}$$

The design matrices $X \in \mathbb{R}^{|\mu|}$ and $W \in \mathbb{R}^{|\Omega|}$ do not necessarily contain distinct explanatory variables.

The outcome regression for the subsample of observed groups needs to condition on $X_G$ and the sample selection rule. It is given by

$$E(Y_G \mid X_G, \text{ selection rule}) \;=\; X_G \beta + \; E(\varepsilon_G \mid \eta_G > \underline{V_G} - W_G \alpha) \tag{11}$$

$$=\; X_G \beta + \; \delta \cdot \lambda(Z_G). \tag{12}$$

In Eqn 11, a consequence of the sample selection rule is that, if groups are selected at random, so that $\varepsilon$ and $\eta$ are independent, the conditional mean of $\varepsilon$ is zero. However, if groups match on observables and unobservables that also affect the outcome, then the conditional mean of $\varepsilon$ is non-zero and $\hat{\beta}$ may be biased. The outcome regression therefore needs to control for the conditional mean of $\varepsilon$, as illustrated in Eqn 12 for the case where the joint distribution of $\varepsilon_G$ and $\eta_G$ is bivariate normal. The last term is the mean of a truncated bivariate normal,

---

[8] In two-group markets with group size $n$, the set of feasible groups includes all $\binom{2n}{n}$ possible $k$-for-$k$ borrower swaps for $k \in \{1, ..., n-1\}$ across the two groups.

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

where $\lambda(\cdot) = \phi(\cdot)/[1 - \Phi(\cdot)]$ is the Inverse Mill's ratio (IMR),[10] as used in the Heckman (1979) selection correction, $Z_G = (\underline{V_G} - W_G\alpha)/\sigma_\eta$, and $\delta$ is a regression coefficient that captures $\rho_{\varepsilon,\eta}/\sigma_\varepsilon$.

## Identification

The outcome equation in the sample selection model requires instrumental variables that determine which groups are formed (instrument relevance) but not the group outcomes (instrument exogeneity). Identifying suitable instruments is challenging when group formation is aimed at optimising outcomes. However, in this model, the interaction in the matching provide this source of identification. No additional instrumental variables are required. Observe that the equilibrium matching is characterised by bounds on the valuations of observed groups. These bounds are a function of the characteristics of other agents in the market. The identification strategy relies on the exclusion restriction that while these characteristics (summarised by the bounds) influence group formation, they do not affect the outcomes of matched groups, which depend solely on the group members. Appendix B provides novel results for point-identification for both matching and outcome equation.

To clarify this concept, consider the example for market 1 in Figure 2a with four agents, denoted by $b$, $c$, $d$ and $e$. The lower bounds summarise agents' characteristics into an instrumental variable. Observe that all agents in the market $(b,c,d,e)$ affect the bounds that affects the matching in the selection equation. For example, for match $bc$, the lower bound $\underline{V_{bc}}$ is affected by the maximum of the characteristics of the counterfactual matches $bd$, $be$, $cd$ and $ce$. But, note that the group's outcome is only determined by the characteristics of its members (i.e. $b$ and $c$). This is the *exclusion restriction*. The characteristics of other agents (i.e. $d$ and $e$) thus provide the *exogenous variation*. Also observe that the introduction of new agents $a$ and $z$ in market 2 and 3 (Figures 2b and 2b), respectively, changes the matching and thus results in groups exhibiting the same outcome but differing characteristics for exogenous reasons.[11]

---

[10]The IMR is the expected value of those values of the error term that cause the group $G$ to be observed.

[11]Group $cd$ in market 2 has the same outcome as group $de$ in market 1, but different correlation (see Figure 1b). Likewise, group $ez$ in market 3 has the same outcome as group $de$ in market 1, but a correlation of one (as opposed to zero) for exogenous reasons.

*Estimation*

In the estimation, the joint distribution of $\varepsilon_G$ and $\eta_G$ is assumed bivariate normal with mean zero and constant covariance $\delta$.

$$\begin{pmatrix} \varepsilon_G \\ \eta_G \end{pmatrix} \sim N \left( 0, \begin{bmatrix} \sigma_\xi^2 + \delta^2 & \delta \\ \delta & 1 \end{bmatrix} \right) \tag{13}$$

Here, $\varepsilon_G = \delta\eta_G + \xi_G$, where $\xi_G$ is a random error. This specification allows for a linear relationship between the error terms in the selection and outcome equations with covariance $\delta$. The variance of the error term of the outcome equation $\sigma_\varepsilon^2$ is $var(\delta\eta + \xi) = \delta^2 + \sigma_\xi^2$. To normalise the parameter scale, the variance of $\eta$ and $\xi$ is set to 1, which simplifies $\sigma_\varepsilon^2$ to $1 + \delta^2$ in the estimation. If the covariance $\delta$ were zero, the marginal distributions of $\varepsilon_G$ and $\eta_G$ would be independent and the selection problem would vanish.

For the estimation, I use 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 D. The conjugate prior distributions of parameters $\alpha$, $\beta$ and $\delta$ 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 $\alpha$ and $\beta$ 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 $\delta$, 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.

### 3.1.1   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

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

(14)

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

(15)

In the matching equation in Eqn 14, the valuation $V_G$ of group $G$ is the sum over all group members' individual utilities (in Eqn 7) from matching with this group.[12] The group valuation is decreasing in liability payment $q$ and increasing in exposure intensity $\epsilon$ and the coincidence of same exposure types. Eqn 15 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 $\eta_G$. The error term $\xi_G$ captures realised individual or aggregate shocks such as health or market demand effects.

The parameter vectors of matching and outcome equation are $\alpha = (r, q, q\epsilon)$ and $\beta = (r, q, q\epsilon, \delta)$, respectively, where the gross interest rate $r$ is known to be constant at 1.09 in the BAAC lending programme and is therefore fixed at this level. The parameters $q$, $q\epsilon$ and $\delta$ are estimated in the model. The expected signs of the parameters are as given in Eqns 14 and 15. 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 or (ii) confounded by a positive selection bias from endogenous group formation. Controlling for unobservable group valuation $\eta_G$ allows me to estimate the treatment effect on repayment net of selection bias. The extent and sign of the selection bias are captured by parameter $\delta$.

---

[12]Note that the group valuation $V_G$ does not contain borrowers' expected returns $E$, because these are constant in the theoretical model and therefore not identified in the empirical model. Further, in the empirical application, we only observe a sample of group members. The terms in the summation signs are therefore normalised to reflect the two-group markets from the theory section. The first term is divided by half the group size $|G|/2$ and the two following terms are divided by $\binom{|G|}{2}$, the number of all possible ways to draw pairs from groups of size $|G|$.

## 3.2   Selection and aggregate effect

In a second step, I test for the adverse selection effect of preventing matching on risk exposure. This effect is estimated in a counterfactual analysis using the coefficient estimates from the sample selection model as parameters.

### 3.2.1   Moral hazard

First reconsider the identification of the moral hazard effects in the previous subsection. Moral hazard manifests in several forms, including ex-ante issues in project selection, as well as ex-post challenges. In each case, increased correlation among borrowers promotes behaviours during and after the loan period that diminish the likelihood of repayment. To attribute differences in repayment outcomes specifically to moral hazard, it is essential to compare two groups that are similar in both observable and unobservable characteristics before receiving loans but are exposed to different levels of correlation. Any observed disparities in repayment between these groups can then be ascribed to the effects of moral hazard. The estimation strategy of the sample selection model in Subsection 3.1 replicates the following ideal experiment with standard cross-sectional survey data.

*Ideal Experiment 1: Treatment effect, net of sample selection bias*

1. Announce in each village that borrower groups will be assigned 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, as announced.

3. Obtain the parameter estimates of randomly and endogenously formed groups. Call the first estimates the *treatment effect* of project covariation and the difference between the two groups the *sample selection bias*.

### 3.2.2   Adverse selection

Adverse selection arises when variations in credit regimes lead to different pools of individuals applying for loans, which can significantly affect repayment. To isolate the adverse selection effect, we require an experimental setup where the credit regime presented to applicants differs between groups, but all other factors remain constant. Crucially, the actual loan terms under which borrowers receive their loans and their awareness of these terms after selection must be identical

across groups. This design ensures that any observed differences in repayment rates can be attributed solely to adverse selection – that is, to the changes in the composition of the applicant pool. This can be thought of as replicating the following ideal experiment.

*Ideal Experiment 2: Selection 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 *selection effect* of matching on risk exposure.

In the second step, individual loans are given out at the end, such that endogenously determined group characteristics do not affect payment and, importantly, we prevent attrition as, all else equal, borrowers prefer individual liability over joint liability.

### 3.2.3 Aggregate effect

To estimate the size of the aggregate effect, I build on the ideal experiment for the selection effect but allow endogenously determined group characteristics, such as project correlation, to affect the payment, i.e. borrowers take joint-liability loans. In particular, I work with the full sample of borrowers in the 2000 BAAC data and run a counterfactual analysis to see how many and what sorts of groups will borrow at the current contract terms under matching regimes (i) and (ii) in the Ideal Experiment 2 above. The participation decision compares utility based on parameter estimates in Eqn 14 to the outside option of agricultural wage labour. The characteristics of the self-selected groups are then used to predict the expected repayment using the parameter estimates from Eqn 15 under both regimes. The protocol for the counterfactual analysis is described in Appendix E.

### 3.2.4 Testable effects and links to theory

The Ghatak (1999) adverse selection model makes two testable predictions for our data. First, it predicts that a reduction in project correlation – e.g. by preventing

the matching on risk exposure – draws existing, safe borrower groups out of the market. Second, this negative effect should be more than offset by the positive effect of the bank securing joint-liability payments more often, as fewer projects fail concurrently.

# 4    Empirical results

The empirical strategy in Section 3 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, 2023b), the corresponding vignette (Klein, 2023a) and in Appendix G. The empirical robustness of the results is examined in Appendix F.

## 4.1    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) |
|---|---|---|
| *Dependent variable* | | |
| - repayment outcome [a] | BAAC *never* raised interest rates as a penalty for late repayment | 0.46 (0.50) |
| *Risk type* | | |
| - risk [b] | Group members' project success prob. | 0.70 (0.07) |
| - risk interaction [b] | Two-way interactions of success prob. | 0.21 (0.03) |
| *Project covariation* | | |
| - same worst year [b] | Measure of coincidence of economically bad years across group members | 0.57 (0.37) |
| *Monitoring cost* | | |
| - same occupation [b] | Measure of occupational homogeneity within group | 0.18 (0.18) |
| *Controls* | | |
| - ln(group age) [c] | Log of number of years group had existed | 4.31 (1.01) |
| - loan size [a] | Average loan size borrowed by the group (thousand Thai baht, currency value in 2000) | 1.59 (1.01) |

[a] from 2000 BAAC group-leader survey

[b] from 2000 BAAC group survey

[c] random regression imputation based on 1997 and 2000 BAAC surveys (see Appendix G)

## 4.2 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 2.3. The average *risk type*, *project covariation* and *monitoring costs* 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 *same
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 mea-
sure establishes a direct link with the different exposure types in Ahlin (2020) 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.

  *Monitoring costs*: The costs of monitoring are measured by the variable *same
occupation*, which gives the borrower-wise coincidence (interaction) of a borrower's
occupational status. Specifically, at the individual level, the variable gives a vector
of the proportion of borrower $i$'s household income accounted for by each of ten
categories.

## 4.3   Treatment effect

The first Probit model in Table 3 reports 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 dummies to control for between-village het-
erogeneity. The resulting positive coefficient suggests that a high level of project
covariation is associated with less arrears. This positive repayment effect can be
explained by either correlation mitigating moral hazard for extremely risk-averse
borrowers (see the Stiglitz model, Proposition 2.2) or the endogenous matching
that biases the coefficient estimate $\widehat{q\epsilon}$ upwards because it picks up the effect of the
omitted risk-type variable. 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 selection bias 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 15).

### 4.3.1   Matching on observables

The above switch in sign implies a positive correlation between risk type and expo-
sure type, which results from endogenous matching on both covariates. Matching

Table 3: Probit and sample selection models with village dummies

| *S.E. in parentheses; one-sided significance at 0.1, 1, 5, 10% denoted by \*\*\*, \*\*, \*, and .* | | | | |
|---|---|---|---|---|
| | **Probit models** | | **Sample selection models** | |
| | (1) | (2) | (3) | (4) |
| **Outcome equation** | | | | |
| *Dependent variable: repayment outcome = 1 if the BAAC has never raised interest* | | | | |
| *as a penalty for late repayment; 0 otherwise.* | | | | |
| *Risk and covariation* | | | | |
| - risk $r$ | – | +1.09 | +1.09 | +1.09 |
| - risk interaction $q$ | – | 0.238 (1.607) | 1.700 (1.713) | 0.028 (2.248) |
| - same worst year $q\epsilon$ | 0.170 (0.289) | -0.015 (0.219) | -0.548 (0.249)* | -1.752 (0.835)* |
| *Monitoring* | | | | |
| - same occupation $1/c$ | – | – | – | -0.455 (0.543) |
| - occ. × worst year $q\epsilon/c$ | – | – | – | 5.983 (3.968). |
| *Controls* | | | | |
| - ln(group age) | -0.040 (0.054) | -0.116 (0.161) | -0.378 (0.117)** | -0.355 (0.115)** |
| - loan size | – | 0.263 (0.421) | 0.911 (0.367)** | 1.109 (0.420)** |
| - loan size sqrd | – | -0.050 (0.088) | -0.180 (0.081)* | -0.210 (0.091)* |
| - village dummies | YES | YES | YES | YES |
| Observations | 68 | 68 | 68 | 68 |
| **Matching equation** | | | | |
| *Dependent variable: group observability indicator = 1 if group is observed; 0 otherwise.* | | | | |
| *Risk and covariation* | | | | |
| - risk $r$ | – | – | -1.09 | -1.09 |
| - risk interaction $q$ | – | – | -1.151 (0.963) | -1.215 (0.985) |
| - same worst year $q\epsilon$ | – | – | 0.290 (0.143)* | 0.285 (0.128)* |
| *Monitoring* | | | | |
| - same occupation $1/c$ | – | – | – | -0.038 (0.476) |
| Observations | – | – | 5,342 | 5,342 |
| **Variance** | | | | |
| Covariance $\delta$ | – | – | 0.453 (0.176)** | 0.513 (0.172)** |

Karlson *et al.* (2012) one-sided test for difference of models (2) and (3), p-value 0.048.
5,284 counterfactual groups and 58 factual groups.

on observables is formally tested in the matching equation of the sample selection 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 drawn for equilibrium and non-equilibrium groups using the Gibbs sampler presented in Subsection 3.1. Turning to the results, the signs of the marginal effects[13] are consistent with the predictions from the theory in Eqn

---

[13]The 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 $Y \sim N(a\mu, a\sigma)$.

14. The negative sign on the risk-type variable means that borrowers value group members with safer projects.[14] While this effect is non-significant, the positive sign on exposure type is significant at the 5%-level and indicates, in line with the theory, that borrowers value peers of the same exposure type. It suggests that exposure type plays an even more significant role in group formation than risk type, which has been the primary focus of the literature.

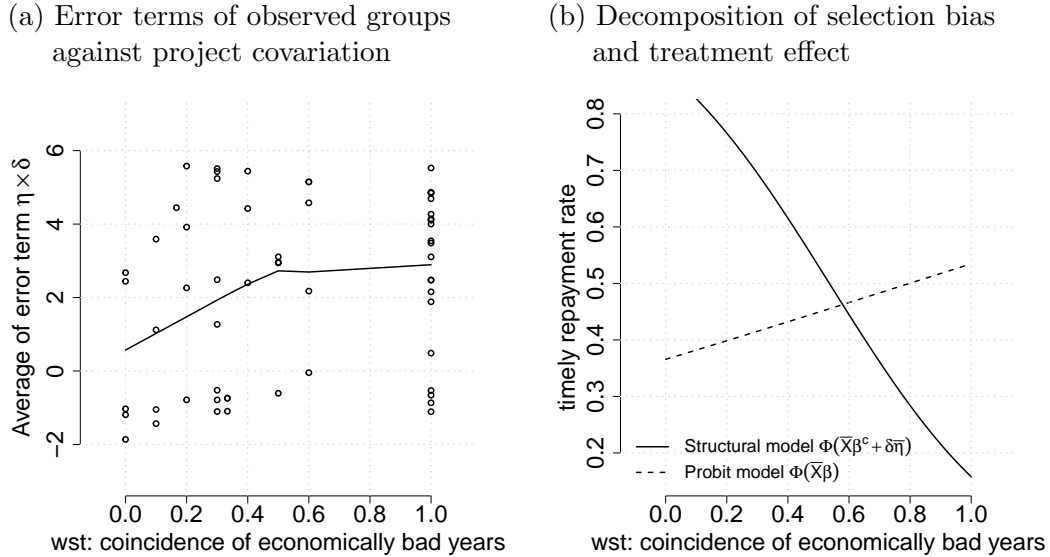### 4.3.2 Matching on unobservables

If matching is also on unobservables that affect group repayment – such as borrowers' local information on risk types – then $\widehat{q\epsilon}$ is still biased upwards in the second Probit model. To correct for this bias, the sample selection 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 $\eta$ 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 $\hat{\delta} = 0.453$, which is equivalent to a correlation of $+0.38$ ($= \frac{\sigma_{\varepsilon,\eta}}{\sigma_\varepsilon \sigma_\eta} = \frac{0.453}{(1+0.453^2)\cdot 1}$). A direct comparison between the second Probit model and the selection-corrected model yields an upwards bias in the Probit model of $+0.53$ ($= -0.015 - [-0.548]$) for $\widehat{q\epsilon}$ that is significant at the 5%-level. This bias results from the positive correlation of project covariation and unobservables $\eta$ 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 sample selection model, the error term $\eta$ in the matching equation enters the outcome equation as $\delta\bar{\eta} > 0$. The omission of this selection-correction term in the Probit regression leads to a positive correlation between project covariation and the error term $\varepsilon$, because $\varepsilon$ is proportional to $\delta\eta$ (i.e. $\varepsilon = \delta\eta + \xi$, where $\xi$ is a random error). Matching on both observables and unobservables thus explains the selection bias in the second Probit model.

### 4.3.3 Decomposition of sample selection bias and treatment effect

Figure 3b illustrates the decomposition of selection bias and treatment effect. The decomposition is done by comparing the estimated regression lines for the first Probit model with the outcome equation of the sample selection model. The models are evaluated conditional on the value of project covariation on the horizontal axis

---

[14]The negative sign on the coefficient results in a positive cross-partial derivative with respect to agents' risk types in Eqn 14.

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



(a) Error terms of observed groups against project covariation

(b) Decomposition of selection bias and treatment effect

with all other variables at their means. The solid regression line of the sample selection model gives the expected repayment – conditional on project covariation – 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 selection bias. To emphasise, if borrowers were assigned at random, as in Ideal Experiment 1 in Section 3, 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. 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 $(= 0.170 \cdot 0.37 = \widehat{q\epsilon}_{probit} \cdot \hat{\sigma})$. This improvement follows from two opposing effects. First, from the sample selection model, we find a significantly negative *treatment effect* of $-20$ percentage points $(-0.548 \cdot 0.37 = \widehat{q\epsilon}_{select} \cdot \hat{\sigma})$ 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 sample selection models we find an even larger but positive *selection bias* of $+27$ percentage points $([\widehat{q\epsilon}_{probit} - \widehat{q\epsilon}_{select}] \cdot \hat{\sigma})$. This

is because the highly correlated groups have unobservables that make them +27 percentage points safer.

### 4.3.4 Decomposition of ex-ante and ex-post moral hazard

In moral hazard models, the idea is that heightened correlation among borrowers promotes behaviour during and after the loan period that diminish the likelihood of repayment. In the theory part, we discuss two forms moral-hazard can take: ex-ante moral hazard in project choice (Stiglitz, 1990) and the ex-post moral hazard problem (Armendáriz, 1999). From our analysis, it is difficult to establish which models are at work. The ex-post model predicts that project covariation's marginal effect on repayment is inversely related to monitoring costs (Corollary 2.2). This prediction is tested by including a monitoring ease variable, $1/c$, in our second sample selection model in Table 3. Although the main effect of monitoring ease is not significant, its interaction with project covariation reaches significance at the 10% level, consistent with the hypothesis that lower monitoring costs enhance the effect of project covariation on repayment. Additionally, the project covariation coefficient, $\widehat{q\epsilon}$, remains significant, indicating that both ex-ante and ex-post moral hazard may be at play in shaping repayment behaviour.

## 4.4 Selection and aggregate effect

For the treatment effect, the empirical model does not allow for an outside option, leaving in or excluding some potential borrower groups. The adverse selection 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 2.2, which predicts a negative repayment effect. To estimate the selection and aggregate effect, I build on the *Ideal Experiment 2* for the selection effect but allow endogenously determined group characteristics, such as project correlation, to affect the payment, i.e. borrowers take joint-liability loans. The participation decision compares expected project returns to the outside option of agricultural wage labour. The protocol for the counterfactual analysis is described in Appendix E.

Table 4 presents the results of the counterfactual analysis for (1) matching on risk type only versus (2) matching on both risk and exposure type. The first row gives the number of groups whose utility from taking a loan exceeds the outside option of agricultural wage labour. Contrary to the predictions of the theory, prohibiting

Table 4: Counterfactual analysis of expected repayment under different matching regimes. Regime (1) is matching on risk type only. Regime (2) is the status quo regime of matching on both risk and exposure type.

| Counterfactuals based on 250 individuals in 29 two-group markets. | | |
|---|---|---|
| *Matching process:* | matching on $p$ only (1) | matching on $p$ and $s$ (2) |
| *Participation* | | |
| No. of groups | 44 | 44 |
| *Group characteristics* | | |
| risk | 0.71 | 0.71 |
| same worst year | 0.57 | 0.70 |
| *Predicted repayment* | | |
| $\widehat{Y}$ | 0.57 | 0.47 |
| 99% CI [a] | (0.52, 0.62) | (0.42, 0.52) |

[a] Confidence intervals based on endpoint transformation.

matching on exposure type (diversification) does not draw more groups out of the programme. In both regimes, the number of participating groups is 44. This can be explained by the small loan sizes in microcredit relative to the expected returns from entrepreneurship. The expected project returns (45,957 Thai baht) for a single borrower exceed the outside option of agricultural wage labour (27,089 baht) by 18,867 baht, which is more than ten times the mean loan size of 1,675 baht.

While diversification does not draw out more borrowing groups, these groups have considerably lower project correlation, as indicated by the average of the share of group members with same worst year falling from 70% to 57% (in row 3). The predicted probability of timely repayment of 0.57 for these groups is consequently higher than when matching on both risk and exposure type (0.47). This is because under lower project correlation, the bank receives more joint-liability payments, consistent with the model predictions from Ghatak (1999). The effect is statistically significant at the 1%-level and of economic importance. Diversification results in a 10 percent increase in timely repayment. A benevolent lender, operating under a zero-profit constraint, would pass this on to the borrowers in the form of a reduction in the interest rate. The results suggest that lenders would benefit from preventing the grouping together of borrowers exposed to similar income shocks.

## 4.5   Discussion

Table 5 compares our findings with those of Ahlin and Townsend (2007), who used a dataset that did not permit analysis of matching into groups. In our theoretical analysis, introducing the negative effect described by Ghatak (2000) reverses the signs in both moral hazard and adverse selection models (Panels A and B). Theory also predicts a positive sample selection bias due to selection on unobservables (Panel C). The discrete choice models (Logit and Probit) report the combined effect of treatment and sample selection bias (Subtotal of A+C). Our sample selection model disentangles these effects, revealing a negative treatment effect and an even larger but positive sample selection bias. The counterfactual analysis in the previous subsection shows that the aggregate effect (Subtotal of A+B) is strictly negative.

Table 5: Summary of theoretical and empirical results

Upward arrows indicate a positive effect of project correlation on expected repayment.

|  | Ahlin/Townsend | | This paper | | | |
|  | Theory | Model | Theory | Models | | |
|  |  |  |  |  | Sample | Counter- |
|  |  | Logit |  | Probit | selection | factual |
| A. Moral hazard |  |  |  |  |  |  |
| - Ex-ante | ↑ |  | ↓[a)] |  | ↓ |  |
| - Ex-post | – |  | ↑ |  | ↑ |  |
| B. Adverse selection | ↑ |  | ↓[b)] |  |  |  |
| C. Selection bias |  |  | ↑ |  | ↑ |  |
| *Subtotal (A+B)* |  |  |  |  |  | ↓ |
| *Subtotal (A+C)* |  | ↑ |  | ↑ |  |  |

[a)] Strictly negative for low risk aversion or low liability payment.
[b)] Strictly negative for low marginal risk types or low liability payment.

These results suggest that banks should discourage the matching of borrowers exposed to similar income shocks. However, policy recommendations must consider whether such restrictions might also prevent borrowers from matching on desirable dimensions, such as social connections. If borrowers tend to match with those they know best, project covariation is naturally linked to social connectedness, since friends or relatives often share income sources and face similar income shocks. Consequently, endogenous matching may produce groups with both correlated returns and social ties.

Regarding optimal market design, three cases emerge. First, if the project correlation measure fully captures social connectedness, group diversification can be implemented by restricting the grouping of relatives, as suggested in the Grameen Replication Guidelines (Alam and Getubig, 2010). The other two cases arise when social connectedness is partially captured in the error term. The implications then depend on the effect of social connectedness on repayment. If social connectedness improves repayment, promoting diversification may have no effect or even a negative impact on repayment. Conversely, if social connections adversely affect repayment, there is a clear case for diversifying groups.

The theoretical and empirical literature offers no consensus on the effect of social connections on repayment. In the context of Thai villages examined in this paper, Ahlin and Townsend (2007) find that cooperative behaviour in groups negatively affects repayment. This aligns with the models of Banerjee *et al.* (1994) and Besley and Coate (1995), which predict that cooperation can prevent a group from exerting repayment pressure on its members. Therefore, the evidence most relevant to our context suggests a positive repayment effect from diversification.

# 5   Conclusion

I analyse the optimal design of rules for group formation in matching markets with an application to joint-liability lending. 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 (treatment effect) and who participates in the market (adverse selection 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 selection bias, I develop a generalised Heckman selection 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 the bias that results from sorting. 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 selection bias even exceeds the negative treatment effect of borrowers' correlated returns on repayment. 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., IRLENBUSCH, B. and RENNER, E. (2006). Group size and social ties in microfinance institutions. *Economic Inquiry*, **44** (4), 614–628.

AGARWAL, N. (2015). An empirical model of the medical match. *American Economic Review*, **105** (7), 1939–1978.

AHLIN, C. (2020). Group lending, matching patterns, and the mystery of microcredit: Evidence from thailand. *Quantitative Economics*, **11** (2), 713–759.

— 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.

ALTINOK, A. (2023). Group lending, sorting, and risk sharing. *Games and Economic Behavior*, **140**, 456–480.

ARMENDÁRIZ, B. (1999). On the design of a credit agreement with peer monitoring. *Journal of Development Economics*, **60** (1), 79–104.

ATTANASIO, O., AUGSBURG, B., DE HAAS, R., FITZSIMONS, E. and HARMGART, H. (2015). The impacts of microfinance: Evidence from joint-liability lending in Mongolia. *American Economic Journal: Applied Economics*, **7** (1), 90–122.

AUMANN, R. and KURZ, M. (1977). Power and taxes. *Econometrica*, **45** (5), 1137–1161.

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.

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.

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.

DE PAULA, A. (2013). Econometric analysis of games with multiple equilibria. *Annual Review of Economics*, **5** (1), 107–131.

DE QUIDT, J., FETZER, T. and GHATAK, M. (2018). Market structure and borrower welfare in microfinance. *The Economic Journal*, **128** (610), 1019–1046.

DUR, U., PATHAK, P. A., SONG, F. and SÖNMEZ, T. (2022). Deduction dilemmas: The taiwan assignment mechanism. *American Economic Journal: Microeconomics*, **14** (1), 164–185.

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.

FERDOWSIAN, A., NIEDERLE, M. and YARIV, L. (2020). *Decentralized matching with aligned preferences*. Tech. rep., Working paper, Department of Economics, Princeton University.[1225].

FOX, J. T. (2010). Identification in matching games. *Quantitative Economics*, **1** (2), 203–254.

— (2018). Estimating matching games with transfers. *Quantitative Economics*, **9** (1), 1–38.

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.

— (1990). Varieties of selection bias. *The American Economic Review*, **80** (2), 313–318.

HERMES, N. and LENSINK, R. (2007). The empirics of microfinance: What do we know? *The Economic Journal*, **117** (517), F1–F10.

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.

KLEIN, T. (2023a). *Analysis of stable matchings in R: Package matchingMarkets*. Vignette to R package matchingMarkets, The Comprehensive R Archive Network.

— (2023b). *matchingMarkets: Analysis of stable matchings*. R package version 1.0-4, The Comprehensive R Archive Network.

KLONNER, S. (2006). *Risky Loans and the Emergence of Rotating Savings and Credit Associations*. Working paper, Cornell University.

MANSKI, C. F. (1985). Semiparametric analysis of discrete response: Asymptotic properties of the maximum score estimator. *Journal of econometrics*, **27** (3), 313–333.

— (1988). Identification of binary response models. *Journal of the American statistical Association*, **83** (403), 729–738.

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.

REED, L. (2015). *Mapping pathways out of poverty: The state of the Microcredit Summit Campaign Report 2015*. Report, Microcredit Summit Campaign.

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.

SØRENSEN, M. (2011). *Identification of multi-index sample selection models*. Working paper, Columbia Business School.

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.

SUTTON, J. (1986). Non-cooperative bargaining theory: An introduction. *The Review of Economic Studies*, **53** (5), 709–724.

TALAMAS, E. (2020). *Nash bargaining with endogenous outside options*. Job Market Paper, IESE Business School.

TOWNSEND, R. (1997a). *Townsend Thai Project Initial Bank for Agriculture and Agricultural Cooperatives (BAAC) Survey, 1997*. Available at http://hdl.handle.net/1902.1/10676, Murray Research Archive.

— (1997b). *Townsend Thai Project Initial Household Survey, 1997*. Available at http://hdl.handle.net/1902.1/10672, Murray Research Archive.

— (2000a). *Townsend Thai Project Bank for Agriculture and Agricultural Cooperatives (BAAC) Annual Resurvey, 2000*. Available at http://hdl.handle.net/1902.1/12057, Murray Research Archive.

— (2000b). *Townsend Thai Project Household Annual Resurvey, 2000*. Available at http://hdl.handle.net/1902.1/10935, 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 2.1.* 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 $\epsilon$.

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

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 4, the expected repayment $\tilde{y}$ of this borrower pool can be written as

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

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

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

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

From Eqn 4 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 distribution[15] yields

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

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

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

---

[15]This implies that $\tilde{p} = \frac{1}{2}(\hat{p} + \underline{p}) = \frac{1}{2}\hat{p}$, $\tilde{\sigma}_p^2 = \frac{1}{12}(\hat{p} - \underline{p})^2 = \frac{1}{12}\hat{p}^2$, $g(\hat{p}) = 1/[1 - \underline{p}] = 1$, 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.[16] In this case, we have $\partial \tilde{y}/\partial \epsilon < 0 \Leftrightarrow q/r < 3/5$. □

*Proof of Corollary 2.1.* The proof of this corollary follows directly from Eqn A4 in the proof of Proposition 2.1. 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

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

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

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

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 (Gangopadhyay *et al.*, 2005). □

*Proof of Proposition 2.2.* The starting point of the proof are two identical, hazardous projects $L$ and $M$ between which borrowers are indifferent.

$$V_{L-M} \quad = \quad V_L - V_M \tag{A10}$$

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

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

Now consider an increase in $\epsilon$ 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?

---

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

Taking the total differential with respect to $\epsilon$ for both projects and allowing a simultaneous change in $p$ and $y$ for the first project yields:

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

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 $dy/dp$.

$$
\begin{aligned}
dp/d\epsilon &= \left\{ -(U_H' - U_{Hq}')dy \right\} / \left\{ (U_H - U_{Hq})2p_H + U_{Hq} + [(U_H' - U_{Hq}')2p_H \right. \\
&\quad \left. + U_{Hq}']dy + [(U_H' - U_{Hq}')(p_H^2 + \epsilon) + U_{Hq}'p_H]\frac{dy}{dp} \right\}.
\end{aligned}
\tag{A13}
$$

The expected repayment to the bank is

$$
Y = r \cdot p_H + q \cdot [p_H(1 - p_H) - \epsilon].
\tag{A14}
$$

Taking the total differential w.r.t. $p$ and $\epsilon$ yields

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

$$
\frac{dY}{d\epsilon} = (r + q(1 - 2p_H)) \cdot \frac{dp}{d\epsilon} - q.
\tag{A16}
$$

Substituting $dp/d\epsilon$ from Eqn A13 above into Eqn A16 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. □

*Proof of Lemma 2.1.*

Let us start with the identification of sub-game perfect equilibria. Agent $i$ with a successful project does not default strategically if the expected pay-off from repayment is larger than that from strategic default, i.e. if

$$
\begin{aligned}
E - r - q(1 - p_j|p_i) &\geq E - \gamma W \tag{A17} \\
\hat{\gamma} &\geq \frac{r + q(1 - p_j|p_i)}{W}, \tag{A18}
\end{aligned}
$$

where, using notation from the previous section and Bayes' rule, the conditional probability, $1 - p_j|p_i$, that agent $i$'s partner $j$ defaults given that agent $i$ is successful can be written as $1 - p_j - \bar{\epsilon}/p_i$. In Eqn A18, the minimum required monitoring effort, $\hat{\gamma}$, is decreasing in project covariation, $\bar{\epsilon}$. □

*Proof of Lemma 2.2.*

Let us now turn back the clock and consider the ex-ante decision by agent $j$ to monitor agent $i$, given the minimum required monitoring effort. Agent $j$ monitors agent $i$ with the required effort, $\hat{\gamma}$, whenever the expected benefit from avoiding liability payment, $q$, for a successful partner who defaults strategically, outweighs the cost of monitoring, $c\hat{\gamma}$. That is, agent $j$ monitors if

$$
q(p_i p_j + \bar{\epsilon}) \geq c\hat{\gamma}. \tag{A19}
$$

Remember that Eqn A19 models an ex-ante monitoring decision and I therefore work with unconditional probabilities. The inequality states that the benefits have to outweigh the costs of monitoring. When investing $c\hat{\gamma}$ in monitoring activities, agent $j$ avoids joint liability payment, $q$, when $i$ is *also* successful but defaults strategically. This event is more likely with correlated project returns. □

*Proof of Proposition 2.3.*

Substitute $\hat{\gamma}$ from Eqn A18 into Eqn A19 to obtain the maximum feasible interest rate

$$r \quad \leq \quad q\bar{\epsilon}\left(\frac{W}{c} + \frac{1}{p_i}\right) - q(1 - p_j) + q\frac{W}{c}p_ip_j \tag{A20}$$

and substitute in the bank's expected repayment in Eqn A21 to obtain

$$
\begin{aligned}
Y_{i,j} &= rp_i + q[p_i(1 - p_j) - \bar{\epsilon}] & \text{(A21)} \\
&= q\bar{\epsilon}\frac{W}{c}p_i + q\bar{\epsilon} - qp_i(1 - p_j) + q\frac{W}{c}p_i^2p_j + qp_i(1 - p_j) - q\bar{\epsilon} & \text{(A22)} \\
&= qp_i\frac{W}{c}(\bar{\epsilon} + p_ip_j). & \text{(A23)}
\end{aligned}
$$

Observe that the marginal effect of project correlation $\bar{\epsilon}$ on repayment is strictly positive

$$\frac{\partial Y_{i,j}}{\partial \bar{\epsilon}} \quad = \quad qp_i\frac{W}{c} > 0. \tag{A24}$$

$\square$

*Proof of Corollary 2.2.*

Observe that the cross-derivative of Eqn A24 is negative: $\frac{\partial Y_{i,j}^2}{\partial \bar{\epsilon} \partial c} = -qp_i\frac{W}{c^2} < 0.$ $\square$

*Proof of Proposition 2.4.* 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 $G \notin \tilde{\mu}$.

$$G \notin \tilde{\mu} \quad \Leftrightarrow \quad V_G < \max_{i \in G} V_{\mu(i)} =: \overline{V_G} \tag{A25}$$

For the *if* direction ($\Rightarrow$) assume for contradiction that $G$ is a *blocking coalition* for $\mu$. 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 $\mu$, 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 $\mu$ 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.

$$G \in \tilde{\mu} \ \Leftrightarrow \ V_G > \max_{G'' \in S} V_{G''} =: \underline{V_G} \tag{A26}$$

Here, $S$ is the set of feasible deviations from $G$, defined as $S(G) := \{G' \in \mathcal{G} | G' \cap G \notin \{\emptyset, G\}, V_{G'} > max_{i \in G' \backslash G} V_{\mu(i)}\}$. 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 $V_{G'} > max_{i \in G' \backslash G} V_{\mu(i)}$. The set of new borrowers are those borrowers in $G'$ that do not belong to the original equilibrium match $G$, i.e. those in $G' \backslash G$.

For the *only if* direction ($\Leftarrow$) assume $\mu$ to be a stable matching with $G \in \mu$. Since $\mu$ 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 $\mu$. It follows from the inequalities in the proposition that no member of $G$ can be part of a blocking coalition. $\square$

*Proof of Proposition B.1.* Let $\alpha_0$ represent the true but unknown parameter vector for the matching model, while $\alpha$ stands for the parameter estimate. We define $\Gamma$ as a set within $R^6$: $\Gamma = \{V \in R^6 : max(V_{ab}, V_{cd}) > max(V_{ac}, V_{ad}, V_{bc}, V_{bd})\}$.

There are $|\Omega| = 3$ possible matchings. We introduce $P_{\mu|W\alpha_0}$ as the probability of observing matching $\mu$, calculated as $P(\eta \in \Gamma - W\alpha_0) = \int 1[\eta \in \Gamma - W\alpha_0]dF(\eta)$. This definition holds for any distribution $F(\eta)$ that adheres to Assumption M1.

We define the partial order, denoted as $\succ_M$ on $R^6$. This order is such that $x \succ_M y$ when

$x = (x_{ab}, x_{cd},\ x_{ac}, x_{bd},\ x_{ad}, x_{bc})$ and

$y = (y_{ab}, y_{cd},\ y_{ac}, y_{bd},\ y_{ad}, y_{bc})$.

The partial order relies on specific conditions, where:

$x_{ab} > y_{ab},\ x_{cd} > y_{cd}$, as well as

$x_{ac} < y_{ac},\ x_{ad} < y_{ad},\ x_{bc} < y_{bc}$ and $x_{bd} < y_{bd}$.

The partial order $\succ_M$ exhibits two crucial properties:

(i) $x \succ_M y \Leftrightarrow x + k \succ_M y + k,\ \forall k \in R^6$ and

(ii) $x \succ_M y \Leftrightarrow x \cdot s \succ_M y \cdot s,\ \forall s > 0$.

We invoke the concept of quantile independence, as described by (Manski, 1988, p. 731). Under this framework, if $W_1\alpha_0 \succ_M W_0\alpha_0$, it implies that $P_{\mu|W_1\alpha_0} > P_{\mu|W_0\alpha_0}$. Specifically, when $W\alpha_0 \succ_M 0$, it leads to $P_{\mu|W\alpha_0} > 1/|\Omega|$ and when $0 \succ_M W\alpha_0$, it results in $P_{\mu|W\alpha_0} < 1/|\Omega|$. For all $\alpha$, we now introduce the set

$$Q_\alpha \ := \ \{W : W\alpha \succ_M 0 \succ_M W\alpha_0\} \cup \{W : W\alpha_0 \succ_M 0 \succ_M W\alpha\} \quad \text{(A27)}$$

Now, we shall discuss two key lemmas that form the basis for identifying $\alpha_0$ relative to $\alpha$.

**Lemma A.1.** *(adapted from Manski, 1988, Proposition 2 and Sørensen, 2011, Lemma 6) Given assumption M1, $\alpha_0$ is identified relative to $\alpha$ when $P(W \in Q_\alpha) > 0$.*

*Proof:* We consider $\alpha$ and $W^0 \in Q_\alpha$. If $P_{\mu|W^0\alpha_0} < 1/|\Omega|$, it implies $P_{\mu|W^0\alpha} > 1/|\Omega|$ for all distributions $F(\eta)$ adhering to M1. Conversely, when $P_{\mu|W^0\alpha_0} > 1/|\Omega|$, it results in $P_{\mu|W^0\alpha} < 1/|\Omega|$ for all such distributions. It is important to note that no distribution function $F(\eta)$ exists that satisfies M1 and leads to $P_{\mu|W^0\alpha} = P_{\mu|W^0\alpha_0}$. This implies that $\alpha$ is identified relative to $\alpha_0$ when $P(W \in Q_\alpha) > 0$.

**Lemma A.2.** *(adapted from Manski, 1985, Lemma 2 and Sørensen, 2011, Lemma 7) Given Assumption M2, $P(W \in Q_\alpha) > 0$ for all $\alpha/||\alpha|| \neq \alpha_0/||\alpha_0||$.*

Here, scale is fixed by selecting a norm $|| \; ||$ on $R^K$, and the proof, as presented, can be readily adpated. The properties in the one-sided matching case do not impact the validity of the proof. For the interested reader, a detailed proof of Lemma A.2 is available from the author upon request. $\qquad \square$

*Proof of Proposition B.2.* Proposition B.1 shows that assumptions M1 and M2 imply A1. Furthermore, we introduce Lemma A.3 below, which shows that M2 implies A4. Hence assumptions A1, A2, A3, and A4 hold, and $\beta_0$ and $F(\varepsilon)$ are identified.

**Lemma A.3.** *When the support $S(W_{ac}, W_{ad}, W_{bc}, W_{bd}, W_{cd})$ does not depend on $X_{ab}$, assumption M2 implies assumption A4.*

*Proof:* Suppose we are given $X_{ab}$, $W_{ab}$ and $\alpha_0$. Without loss of generality, we can assume that the characteristics in $X_{ab}$ are dependent on the characteristics in $W_{ab}$, and $W_{ab}$ corresponds to $W_C$. Under the conditions of M2(1) and the assumption that the support of $(W_{ac}, W_{ad}, W_{bc}, W_{bd}, W_{cd})$ is independent of $X_{ab}$, the variables $(W_{ac}, W_{ad}, W_{bc}, W_{bd}, W_{cd})$ corresponds to the subset $W_I$. These variables provide the necessary independent variation for identification.

We proceed to select an increasing sequence $k^i$, such that the limit as $k$ approaches infinity is infinite. With assumption M2(3), we can select $W_0^i$ and $W_1^i$ from the support of $G_0$ such that $-W_0^i \alpha_0 = k^i$ and $-W_1^i \alpha_0 = -k^i$. To facilitate this, we define $W^i$ as follows: $W^i = (W_{ab}, W_{ac}, W_{ad}, W_{bc}, W_{bd}, W_{cd}) = (W_{ab}, W_0, W_1, W_1, W_1, W_1)$. We establish $A^i$ as $A^i = \Gamma - W^i \alpha_0 = \{V \in R^6 : max(V_{ab} - W_{ab}\alpha_0, V_{cd} + k^i) > max(V_{ac} - k^i, V_{ad} - k^i, V_{bc} - k^i, V_{bd} - k^i)\}$. Crucially, $A^i$ is a subset of $A^{i+1}$ and, most significantly, $P(\eta \in \cup_i A^i) = 1$. This satifies the conditions of assumption A4 and establishes that M2 implies A4. $\qquad \square$

# B  Identification

## B.1  Identification of the one-sided matching model

The one-sided matching model without transfers belongs to the category of discrete choice models. Despite the significance of this model, there is a notable gap in existing research concerning identification. However, it is worth mentioning that Jeremy Fox has analysed the identification (Fox, 2010) and estimation (Fox, 2018) of matching games with transfers, which is a related field. We bridge this gap by deriving identification results based on prior findings related to disrete choice models (Manski, 1985, 1988) and two-sided matching markets (Sørensen, 2011).

We continue with the example introduced in Section 1, involving four agents labelled as $a$, $b$, $c$ and $d$. We will explore two essential distributions:

- $F(\eta)$ represents the distribution of $\eta$, which is a vector composed of error terms for all possible matches of agents, i.e. $\eta = (\eta_{ab}, \eta_{ac}, \eta_{ad}, \eta_{bc}, \eta_{bd}, \eta_{cd}) \in R^6$

- $G(W)$ denotes the distribution of $W$, which is a matrix containing characteristics for all possible matches, i.e. $W = (W_{ab}, W_{ac}, W_{ad}, W_{bc}, W_{bd}, W_{cd}) \in R^{6 \times K}$.

Our identification results are built upon two key distributional assumptions:

**M1:** Let $F(\eta) = F_0 \times F_0 \times F_0 \times F_0 \times F_0 \times F_0$, where $F_0$ represents an absolute continuous distribution, and its support extends across the real line $R^1$.

**M2:** This assumption has multiple elements:

(1) Let $G(W) = G_0 \times G_0 \times G_0 \times G_0 \times G_0 \times G_0$, with $G_0$ representing an absolute continuous distribution on $R^K$.

(2) With probability 1, the support of $G_0$ is not contained within any proper subspace of $R^K$.

(3) For at least one value of $k$ in the range $(1, 2, ..., K)$, we observe that, for almost every value of $\tilde{W} = (W^1, ..., W^{k-1}, W^{k+1}, ..., W^K)$, the distribution of $W^k$ conditional on $\tilde{W}$ exhibits positive Lebesgue density under $G_0$.

Assumption M1 is a prominent distributional assumption. In the context of the one-sided matching model, it brings forth three crucial implications. (1) It implies

that the six elements of $\eta$ are independent and identically distributed. (2) It ensures that $P(\mu|W\alpha_0 = 0) = 1/|\Omega|$, i.e. the probability of picking a matching $\mu$ at random is equal to the inverse of the number of possible matchings. In the context of the example introduced in Section 1, we have $|\Omega| = 3$. (3) It guarantees that $0 < P_{\mu|W\alpha_0} < 1 \; \forall \; W$. Assumption M2 corresponds to Assumption 2 in Manski (1985) and specifies the required variation in $W$.

**Proposition B.1.** *Under assumptions M1 and M2, the parameter vector $\alpha_0$ is identified up to scale.*

*Proof:* The proof to this proposition extends the results in Manski (1985, Lemma 2), Manski (1988, Proposition 2) and Sørensen (2011, Lemma 4), and is provided in Appendix A.

## B.2   Identification of the outcome equation

The economic analysis of selection in the context of one-sided matching employs a framework known as the multi-index sample selection model. This model, introduced in Sørensen (2011), not only characterizes selection but also studies identification. Identification is a key concern in the estimation of such models, as it necessitates exogenous variation. In this model, the variation is provided through the characteristics of agents who share the same market but do not belong to the same group.

In the context of this model, it is pivotal to discern between the elements of the variable $W$. Specifically, some elements of $W$ coincide with elements of another variable, denoted as $X$. Thus, the matrix $W$ can be divided into two subsets: common elements, referred to as $W_C$, and independent elements, termed $W_I$. The latter category, $W_I$, encompasses those elements of $W$ that exhibit variation when $X$ is held constant, while $W_C$ comprises elements of $W$ that remain fixed when $X$ is kept constant.

To facilitate the discussion, we introduce the following set

$$\tilde{S}(X, W_C) \; := \; \{\Gamma - W\alpha_0 : (W_C, W_I) \in S(W|X)\}. \qquad (A28)$$

$\tilde{S}(X, W_C)$ is the collection of all possible sets defined as $\Gamma - W\alpha_0$, where $(W_C, W_I)$ belongs to $S(W|X)$, the support of $W$ conditional on $X$. Here, $X$ and $W_C$ are held constant, and $\Gamma - W\alpha$ varies solely based on the elements of $W_I$. Essentially, this set encapsulates the variation in $\Gamma - W\alpha$ that is induced by changes in the elements of $W_I$. It is akin to an exclusion restriction where variables in $W_I$ are excluded

from the outcome equation. The selection equation postulates that $Y$ is observed when the error terms fall within a specific set, i.e. $\eta \in \Gamma - W\alpha$. $\tilde{S}(X, W_C)$ embodies the variation in this set, brought about by alterations in $W_I$. This independent variation serves as the crux for identifying the outcome equation in this model. A crucial condition for the identification of $Y$ is that the set $\tilde{S}(X, W_C)$ exhibits sufficient variation

To clarify this concept, consider the initial example in Section 1 involving four agents, denoted as $a$, $b$, $c$ and $d$. For the outcome equation of group $ab$, the characteristics are denoted by $X = (X_{ab})$. The elements that enter the matching equations encompass $W = (X_{ab}, X_{cd}, X_{ac}, X_{ad}, X_{bc}, X_{bd})$. The independent components of $W$ are represented as $W_I = (X_{cd}, X_{ac}, X_{ad}, X_{bc}, X_{bd})$. For all conceivable values of $X_{ab}$, the set $\tilde{S}(X, W_C)$ comprises subsets, $S \subset R^4$, such that when $X\alpha + \eta$ falls within $S$, it satisfies the condition $max(X_{ab}\alpha + \eta_{ab}, X_{cd}\alpha + \eta_{cd}) > max(X_{ac}\alpha + \eta_{ac}, X_{ad}\alpha + \eta_{ad}, X_{bc}\alpha + \eta_{bc}, X_{bd}\alpha + \eta_{bd})$. This condition depends on the choice of $X_{cd}$, $X_{ac}$, $X_{ad}$, $X_{bc}$ and $X_{bd}$.

Sørensen (2011) outlines four assumptions, namely A1, A2, A3, and A4, which underpin the identification of the parameters $\beta_0$ and the distribution $F(\varepsilon)$ for multi-index sample selection models. These assumptions are as follows:

**A1:** $\Gamma(W)$ is known.

**A2:** This assumption encompasses two elements: (1) The distribution of $(\varepsilon, \eta)$ is independent of $W$, and (2) the distribution $F(\varepsilon)$ is absolute continuous with support equal to the real line.

**A3:** $F(\varepsilon)$ either has (1) mean zero, or (2) median zero.

**A4:** For all $X$ and $W_C$, there exists a sequence of sets $A_i \in \tilde{S}(X, W_C)$, such that $A_i \subset A_{i+1}$, $P(\eta \in A_i) > 0$, and $P(\eta \in \cup_i A_i) = 1$.

Assumption A1 states that the parameters governing the selection process are identified and known. Assumptions A2 and A3 are standard technical assumptions. Assumption A4 states that $\tilde{S}(X, W_C)$ is sufficiently rich in the sense that it contains a given collection of subsets, and this defines the required independent variation in $W_I$.

**Proposition B.2.** *Given assumptions M1, M2, A2 and A3, both $\beta_0$ and $F(\varepsilon)$ are identified.*

*Proof:* The proof to this proposition extends the results in Heckman (1990) and Sørensen (2011, Proposition 4), and is provided in Appendix A.

# C   Bargaining model

The bargaining model at the group formation stage can be formulated as

$$
\begin{aligned}
\max_{\tilde{c}} \quad & \Pi_{i=1}^{n} u_{i,G}(\tilde{c}_i) - \omega_i \\
\text{s.t.} \quad & \sum_{i=1}^{n} \tilde{c}_i && = && x \\
& \omega_i && \leq && u_{i,G}(\tilde{c}_i) && \forall i \in G \\
& \tilde{c}_i && \geq && 0 && \forall i \in G,
\end{aligned}
$$

where $\omega_i$ is borrower $i$'s outside option and $x$ is a vector of net group payoffs. To illustrate, for $n = 2$ players, $i$ and $j$, we have $x = (x_{ij}, x_i, x_j, 0)$, where $x_{ij}$ is the payoff when both players' projects succeed, and $x_i$ and $x_j$ are the payoffs when only player $i$ or player $j$ succeeds, respectively. For $n$ players, the length of this power set is $2^n$. Further, $\tilde{c}_i = (\tilde{c}_{i,ij}, \tilde{c}_{i,i}, \tilde{c}_{i,j}, 0)$ is the vector of player $i$'s share of the group payoff $x$. The vector $\tilde{c} = (\tilde{c}_1, ..., \tilde{c}_n)$ is the collection of all borrowers' shares. For a group of two borrowers, the payoffs and success probabilities are

- $x_{ij} = y_i + y_j - 2r$ with $P_{ij} := p_i p_j + \epsilon$, if both players are solvent[17]

- $x_i = y_i - r - q$ with $P_i := p_i(1 - p_j) - \epsilon$, if only $i$ is solvent

- $x_j = y_j - r - q$ with $P_j := (1 - p_i)p_j - \epsilon$, if only $j$ is solvent

Both players aim to maximise their expected utilities, which are given by

- $u_{i,j}(\tilde{c}_i) = P_{ij} u(\tilde{c}_{i,ij}) + P_i u(\tilde{c}_{i,i}) + P_j u(\tilde{c}_{i,j})$

- $u_{j,i}(\tilde{c}_j) = u_{j,i}(x - \tilde{c}_i) = P_{ij} u(x_{ij} - \tilde{c}_{i,ij}) + P_i u(x_i - \tilde{c}_{i,i}) + P_j u(x_j - \tilde{c}_{i,j})$

and sum up to the group valuation $u_{i,j}(\tilde{c}_i) + u_{j,i}(\tilde{c}_j) = V_{ij}$. Using the Lagrange multiplier method and taking first order conditions, we have

$$
\frac{u'(\tilde{c}_{i,ij})}{u'(x_{ij} - \tilde{c}_{i,ij})} = \frac{u'(\tilde{c}_{i,i})}{u'(x_i - \tilde{c}_{i,i})} = \frac{u'(\tilde{c}_{i,j})}{u'(x_j - \tilde{c}_{i,j})} = \frac{u_{i,j}(\tilde{c}_i) - \omega_i}{u_{j,i}(x - \tilde{c}_i) - \omega_j}. \tag{A29}
$$

Rewriting this in terms of the fear of ruin, $\chi$, (Aumann and Kurz, 1977) as follows

$$
\frac{u_{j,i}(x - \tilde{c}_i) - \omega_j}{u'(x_{ij} - \tilde{c}_{i,ij})} = \chi_j(x_{ij} - \tilde{c}_{i,ij}) = \chi_i(\tilde{c}_{i,ij}) = \frac{u_{i,j}(\tilde{c}_i) - \omega_i}{u'(\tilde{c}_{i,ij})}, \tag{A30}
$$

we find that both players have the same fear of ruin, which implies equal sharing and a unique stable matching. By the symmetry property of Nash bargaining, this also applies for groups with mixed exposure types and size larger than two.

---

[17]The probability $P_{ij}$ is obtained as $Pr(y_i = 1, y_j = 1) = \frac{1}{2}(p_i + \tilde{\gamma})(p_j + \tilde{\gamma}) + \frac{1}{2}(p_i - \tilde{\gamma})(p_j - \tilde{\gamma}) = p_i p_j + \tilde{\gamma}^2 = p_i p_j + \epsilon$.

# D    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 nuissance 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 $N\left(X_G\beta + (V_G - W_G\alpha)\delta, 1\right)$ with density

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

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 $V_G$, $\alpha$ and $\delta$ 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 $N\left(W_G\alpha, 1\right)$, truncated from above at $\overline{V_G}$ with density

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

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

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

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, $\sigma_\eta^2$, equals one.[18]

## Conditional posterior distribution of parameters

### Alpha

The coefficient vector $\alpha$ 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 $\alpha$ is $N\left(\hat{\alpha}, \hat{\Sigma}_\alpha\right)$, where

$$\hat{\Sigma}_\alpha = \left[\Sigma_\alpha^{-1} + \sum_{t \in T_2}\left[\sum_{G \in M_t} W_G'W_G + \sum_{G \in \mu_t} \frac{\delta^2}{\sigma_\xi^2} W_G'W_G\right]\right]^{-1} \tag{A31}$$

and

$$\hat{\alpha} = -\hat{\Sigma}_\alpha\left[-\Sigma_\alpha^{-1}\bar{\alpha} + \sum_{t \in T_2}\left[\sum_{G \in M_t} -W_G'V_G\right.\right.$$
$$\left.\left. + \sum_{G \in \mu_t} \frac{\delta}{\sigma_\xi^2} W_G'\left(Y_G^* - X_G\beta - V_G\delta\right)\right]\right]. \tag{A32}$$

The variables $\Sigma_\alpha^{-1}$ 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 $\beta$ is $N\left(\hat{\beta}, \hat{\Sigma}_\beta\right)$, where

$$\hat{\Sigma}_\beta = \left[\Sigma_\beta^{-1} + \sum_{t \in T_1}\sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X_G'X_G + \sum_{t \in T_2}\sum_{G \in \mu_t} \frac{1}{\sigma_\xi^2} X_G'X_G\right]^{-1} \tag{A33}$$

---

[18] $\sigma_\eta^2 = var(\frac{\varepsilon\delta}{\sigma_\xi^2 + \delta^2} + x) = \frac{(\sigma_\xi^2 + \delta^2)\delta^2}{(\sigma_\xi^2 + \delta^2)^2} + \sigma_x^2 = \frac{\delta^2}{(\sigma_\xi^2 + \delta^2)} + \sigma_x^2$. So $\sigma_\eta^2 = 1$ iff $\sigma_x^2 = \sigma_\xi^2/(\sigma_\xi^2 + \delta^2)$.

and

$$
\begin{aligned}
\hat{\beta} \; = \; -\hat{\Sigma}_\beta \Bigg[ &-\Sigma_\beta^{-1}\bar{\beta} - \sum_{t\in T_1}\sum_{G\in\mu_t}\frac{1}{\sigma_\xi^2}X'_G Y^*_G \\
&-\sum_{t\in T_2}\sum_{G\in\mu_t}\frac{1}{\sigma_\xi^2}X'_G(Y^*_G - \delta(V_G - W_G\alpha))\Bigg].
\end{aligned}
\tag{A34}
$$

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

$$
\hat{\sigma}_\delta^2 \; = \; \left[ \frac{1}{\sigma_\delta^2} + \sum_{t\in T_2}\sum_{G\in\mu_t}\frac{1}{\sigma_\xi^2}(V_G - W_G\alpha)^2 \right]^{-1}
\tag{A35}
$$

and

$$
\hat{\delta} \; = \; -\hat{\sigma}_\delta^2 \left[ -\frac{\bar{\delta}}{\sigma_\delta^2} - \sum_{t\in T_2}\sum_{G\in\mu_t}\frac{1}{\sigma_\xi^2}(Y^*_G - X_G\beta)(V_G - W_G\alpha) \right].
\tag{A36}
$$

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

# E    Protocol for counterfactual analysis

The counterfactual analysis follows the protocol below and is based on the coefficient estimates.

1. **Group formation.** Obtain the equilibrium groups in the 29 two-group markets when matching is on either of the two regimes below:

    P+S: risk and exposure type (status quo regime)

    P: risk type only (diversification regime).

    The second regime is akin to preventing the matching on exposure type or, equivalently, setting $\epsilon = 0$ in Eqn 14. Equilibrium groups are determined using the group valuation in this equation. and a top-down algorithm (Talamas, 2020).

2. **Selection decision.** Group $G$ takes a loan if the expected project returns of the group $n \cdot E_t$ in market $t$ minus the expected interest and liability payment $l_t \cdot \hat{V}_G$ exceeds the outside option $n \cdot \bar{u}_t$, taken as the average wage rate for agricultural labour.

$$n \cdot E_t + l_t \cdot \hat{V}_G \quad > \quad n \cdot \bar{u}_t, \tag{A37}$$

    where $n$ is the group size, the expected return $E_t$ is obtained from the BAAC group survey, that asks group members for their expected income for the following year. $l_t$ is the average loan size in market $t$ and $\hat{V}_G$ is the predicted group valuation with estimates from Eqn 14.

3. **Expected repayment.** For the remaining groups, predict the expected repayment for regimes $P + S$ and $P$, using the parameter estimates from Eqn 15 in the sample selection model.

# F   Robustness of the results

This appendix examines whether the primary result – the decomposition into a negative treatment effect and a positive selection 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 -0.055, p-value 0.653, meaning that survivorship bias is not an issue.

Third, 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 (2023a) presents Monte Carlo evidence of the robustness of the estimator in small samples, which confirms that the resulting attenuation bias even underrates the selection bias that this paper corrects for.

Furthermore, the surveys of the Townsend Thai project sample up to two groups per village. This number concurs with the average number of groups per village. In the 1997 household survey, 22% households indicate having obtained a BAAC group-guaranteed loan in the past year. With an average of 121 households per village and an average group size of 11.5 borrowers, there are an average of $121 \cdot 0.22/11.5 = 2.3$ groups per village.

Finally, I consider the relevance of the instrumental variables. For the instrument (summarised by the bounds) to be relevant, they need to be a good predictor of the selection decision. We cannot observe this, because we do not estimate a

coefficient on the bounds. Instead, we test whether the inverse Mill's ratio (IMR) is correlated with variables in the outcome equation. If the instrument is not relevant, then the model is only identified by the curvature of the normal density function and the IMR is almost collinear with $X\beta$.

In Heckman (1979), the IMR is $\lambda(-W\alpha)$. When the characteristics in the selection equation, $W$, are the same as those in the outcome equation, $X$, then identification is weak, because it relies exclusively on the non-linearity of the function $\lambda(\cdot)$. This case with instruments is depicted with the black dots in Figure A1, where each dot represents one of the 68 observed groups. In the sample selection model in this article, we also have $W = X$, but the IMR is $\lambda(\underline{V_G} - W\alpha)$ and the lower bounds $\underline{V_G}$ serve the role of an instrumental variable. This case with an instrument is represented using the white dots in Figure A1. For half of the groups, the IMR is zero because the lower bound for these groups is $-\infty$. For the other half of the groups, the equilibrium bounds provide significant independent variation. Statistical test of the null hypothesis of a linear association between IMR and the linear predictor $X\beta$ in Table A1 shows that the null is rejected for the case with an instrument (white dots) but not for the case without an instrument (black dots).

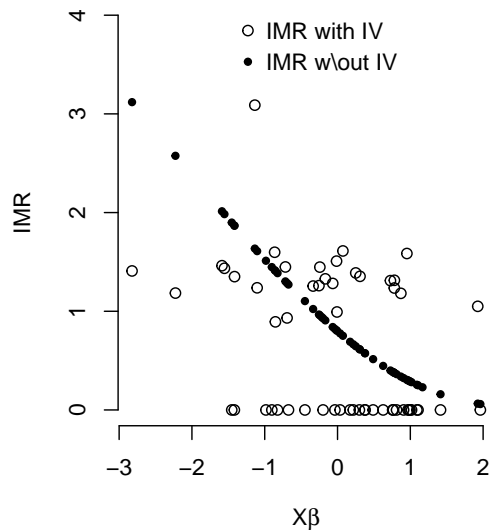Figure A1: Comparison of Inverse Mill's Ratio (IMR) with and without IV



Table A1: Significance tests for the null of perfect correlation between IMR and $X\beta$

|  | with IV | without IV |
|---|---|---|
| $cor_{IMR,X\beta}$ | $-0.128$ | $-0.982$ |
| $cor_{IMR,X\beta} + 1$ | 0.872 | 0.018 |
| t-statistic | 6.580 | 0.704 |
| p-value | 0.000 | 0.242 |
| N | 58 | 58 |

# G     Replication guide[19]

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`.

## G.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.

### G.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.

### G.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 BAAC 1997 and 2000 surveys in an interval regression. While the group age is not observed in

---

[19]This section of the Appendix is not intended for publication.

the BAAC 2000 data, the quasi-panel still allows me to find bounds for a group's age (see Table A2 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 A2: Definition of bounds for interval regression of the missing group_age variable

| Groups from | lower bound | upper bound |
|---|:---:|:---:|
| *BAAC 1997 survey* | $group\_age_{97}$ | $group\_age_{97}$ |
| *BAAC 2000 survey* | | |
| - in villages with single group in '97 | $\max\{group\_hist_{00}, 1\}$ | $max\_age_{97}+3$ |
| - in all other villages | $\max\{group\_hist_{00}, 1\}$ | 2000-1966 |

The results of the interval regression are presented in Table A3 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 A2a and A2b 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 A2. 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.

### G.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`.

### G.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, 2023b). The resulting group-level data is in `data/group-level.RData`.

## G.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 LaTeX 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.
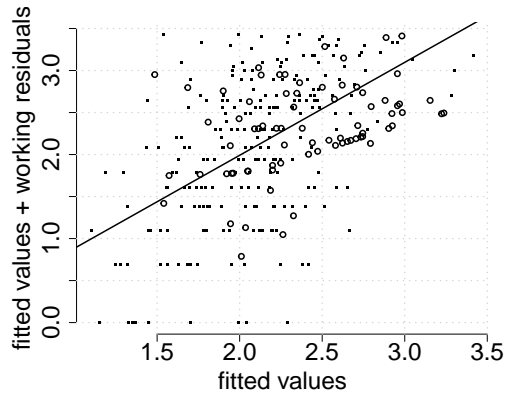
The estimator developed in this paper is implemented in the R package and the source code available on the Comprehensive R Archive Network. To test the functionality of the software implementation in this package, Klein (2023a) provides simulation evidence of the correct implementation of both design matrix generation and estimators.

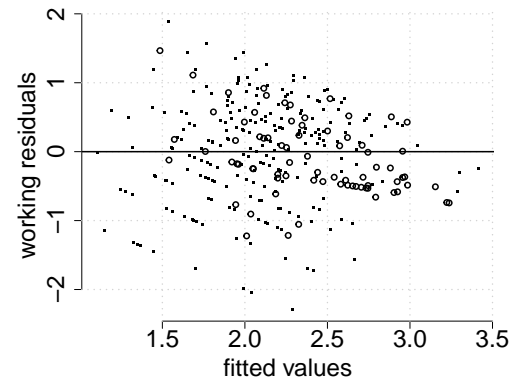Table A3: Interval regression imputation of the missing group_age variable

| | |
|---|---|
| *S.E. in parentheses; significance at 0.1, 1, 5, 10% denoted by \*\*\*, \*\*, \*, and . respectively.* | |
| **Interval regression** | |
| *Dependent variable as defined in Table A2.* | |
| 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, $Pr(> \chi^2_7)$ | 1e–14 |

Figure A2: 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

Download ZEW Discussion Papers:

https://www.zew.de/en/publications/zew-discussion-papers

or see:

https://www.ssrn.com/link/ZEW-Ctr-Euro-Econ-Research.html
https://ideas.repec.org/s/zbw/zewdip.html

//

IMPRINT