

# Credit Allocation Mix in Peer-to-Peer Lending: A Network Study on Bondora

Social Network Analysis - Group 7

Floris Vermeulen      Martijn van Iterson      Niek Fleerakkers  
Patryk Grodek      Samir Sabitli

2025-11-16

## Table of Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Dataset</b>	<b>5</b>
2.1	Pre-Processing & Network Formation . . . . .	6
2.2	Descriptive Statistics & Preliminary Analysis . . . . .	6
<b>3</b>	<b>Research Rationale</b>	<b>8</b>
<b>A</b>	<b>Supplements</b>	<b>10</b>
A.1	Data Preprocessing Steps . . . . .	10
A.2	Distributions of Variables across Samples . . . . .	12
A.3	Source Code - Data Preprocessing . . . . .	15
A.4	Source Code - QAP Linear Regression . . . . .	21
A.5	Source Code - ERGM Network Analysis . . . . .	23
	<b>References</b>	<b>28</b>
<b>B</b>	<b>Technology Statement</b>	<b>30</b>

## List of Figures

1	Network Visualisation . . . . .	7
2	Network Centrality Distribution Plots . . . . .	8
3	Number of Unique Users per Loan Type per Year . . . . .	10
4	Distribution of Loan Amount across Samples . . . . .	12
5	Distribution of Interest Rates across Samples . . . . .	12
6	Distribution of Age across Samples . . . . .	12
7	Distribution of Loan Duration across Samples . . . . .	13
8	Distribution of Loan Purpose across Samples . . . . .	13

9	Distribution of Credit Ratings across Samples . . . . .	13
10	Distribution of Occupation Types across Samples . . . . .	14

## List of Tables

1	Hypotheses related to Study 2 . . . . .	5
2	Descriptive Statistics of Numeric Variables . . . . .	6
3	Basic Network Descriptive Statistics . . . . .	7

# 1 Introduction

With the advent of online peer-to-peer (P2P) lending platforms, the traditional methods of financial intermediation have been usurped by individual choice; with no stringent third parties involved in transactions, individuals have greater power in sourcing credit. This presents new realities for individuals disenfranchised by the traditional financial system (Alistair Milne & Paul Parboteeah, 2017). Traditional financial literature shows that individually, investors exhibit behavioural biases (Agrawal, 2012; Ayal, Hochman, & Zakay, 2011) in how they choose their investments, however, numerous studies demonstrate that these behavioural biases are present among non-professional P2P borrowers. For example, Herzenstein, Dholakia, & Andrews (2011) and Lee & Lee (2012) suggest that users cluster around popular loans, exhibiting ‘herding’ behaviour. Literature focuses predominantly on funding success (Yao, Chen, Wei, Chen, & Yang, 2019), however, Ayal, Bar-Haim, & Ofir (2018) suggests that there is insufficient attention on how borrowers form their loan portfolios. As the authors indicate, this area of research is important in understanding why less professional investors deviate from known models of rational investor behaviour. Our study aims to model and understand how P2P borrowers’ behaviours influence what kinds of loans they acquire. In subsequent discussions, we consider investors as rational participants in credit markets, as described by theory, whereas borrowers refer to the users of P2P platforms.

Investigating behavioural biases, Ayal et al. (2018) find that investors’ familiarity with specific assets can make them seem less risky, making recognisable assets more attractive. However, such assets lead to undiversified portfolios. Evidence with P2P platforms supports the familiarity bias in lenders, as Galak, Small, & Stephen (2010) demonstrate that lenders choose borrowers who are similar across demographic attributes such as gender, occupation, and ethnicity. Based on these notions, we formulate our first study, asking:

**Research Question 1:** *Do borrowers tend to choose similar loan uses based on their individual and loan characteristics?*

This question is answered using linear Quadratic Assignment Procedure (QAP) regressions, as we can represent similarities across loan uses and attributes as matrices that can be used as inputs in the model. Further, we integrate loan use frequencies to understand how the prevalence of specific loan types affects borrowers’ portfolio compositions.

Using data from LendingClub, Serrano-Cinca, Gutiérrez-Nieto, & López-Palacios (2015) study funding success, finding that the reported loan use is a significant factor explaining differences in default rates. Specifically, the authors argue that there is heterogeneity in the credit riskiness of individuals seeking specific loan purposes, indicating that credit rating can explain some variation in borrowers’ choice of loan use. Therefore, we formulate the following hypothesis:

*Hypothesis 1: Borrowers with the same credit ratings choose similar loan uses*

Supporting the familiarity bias, Ravina (2019) studies personal characteristics in P2P markets, finding that loans tend to be awarded to those sharing personal characteristics such as occupation. As Ayal et al. (2018) implies, this can be generalised to suggest that lenders’ loan grants reflect a wider network of similar borrower preferences, which can reflect the homogeneity in loan use observed by Serrano-Cinca et al. (2015). Thus,

*Hypothesis 2: Borrowers sharing the same occupation area tend to share loan uses*

Further, behavioural finance suggests that individuals’ financial decisions are not independent, but rather, shaped by bounded rationality and psychological preferences (Kahneman & Tversky, 1979). In the P2P context, Ayal et al. (2018) shows that these preferences can influence how borrowers choose loan types and how this can overlap with other borrowers’ decisions. Using the Exponential Random Graph Model (ERGM) framework, we can identify how these behaviours manifest as structural regularities in a borrower-to-loan-use network. Therefore, we ask:

**Research Question 2:** *What behavioural mechanisms affect the structure of borrower-to-loan-use connections?*

The aforescribed section outlined that the familiarity bias can lead to similarities across borrower portfolios in terms of their characteristics. However, in P2P settings, this can manifest itself as familiarity *clustering*, wherein borrowers who have previously been granted certain loan types are more likely to repeat their choices, or select other categories that appear familiar to different borrowers. Herzenstein et al. (2011), studying borrower decisions on Prosper.com, identifies a ‘strategic herding’ behaviour, where individuals gravitate towards popular or socially-validated loans. Considering that loan popularity is directly observable through its number of bids, we assert that borrowers linked to one loan purpose are likely to connect with others sharing the same loan use. In network terms, this behaviour results in cyclical substructures where there are overlapping nodes of each of the bipartite partitions. We capture a conservative form of this cycle through the ERGM term `cycle(4)`, formulating:

*Hypothesis 3: Borrowers who share one loan purpose are likely to share another*

Further, according to Amar, Ariely, Ayal, Cryder, & Rick (2011), the notion of Debt Account Aversion can also explain how investors choose debt. Specifically, investors tend to avoid holding multiple debt accounts because per Prospect Theory, borrowers mentally segregate wins and losses to minimise perceived distress from debt (Kahneman & Tversky, 1979). Ayal & Zakay (2009) expands with the Theory of Perceived Diversification, explaining that borrowers’ distress increases with the perceived *distinctiveness* of each debt type, suggesting that in the P2P context, debt mentally differentiated by purpose would also be seen as distressing (Ayal & Zakay, 2009). Therefore, we assert that borrowers try to minimise loan distinctiveness by requesting loans of a single type. This can be captured by the `b1degree(1)` term, which models whether borrowers are only linked to one loan purpose. Therefore, we formulate:

*Hypothesis 4: Borrowers exhibit preferences for minimal debt distinctiveness, selecting one loan use*

Notwithstanding, credit choices relate to non-structural factors. Cooper, Gorbachev, & Luengo-Prado (2023) show that younger borrowers face constraints in the types of loans they can access in typical credit markets. The authors argue that this constraint shapes how these borrowers form loan portfolios, often concentrating around different use cases such as education, personal consumption, or small business loans. Considering the systemic differences in how age affects credit consumption, we assert that younger borrowers have different loan use portfolios. The ERGM term `b1cov("age")` can capture this by assessing how changes in age affect how borrowers form connections with different loan types. Therefore:

*Hypothesis 5: Younger borrowers have different loan use mixes than older borrowers*

Alongside age, the borrower’s gender is known to be a significant factor in influencing credit use decisions. Aliano, Alnabulsi, Cestari, & Ragni (2023) study the role of gender and other factors on

borrowers' probability of default on the Bondora.com, finding a persistent gender effect on different loan uses; women tend to default less on health, home, and business loans. Croson & Gneezy (2009) shows that these differences can be due to risk perception and self-selection effects, where women exhibit greater risk aversion, leading to different borrowing patterns. We expect that these differences in risk tolerance and perceived creditworthiness by lenders affect loan use composition in terms of gender. The term `b1nodematch("gender")` captures this effect by assessing homophily in loan use selection. Our premise is supported by a negative estimate for the term, indicating heterophily. Therefore:

*Hypothesis 6: Gender differences lead to different loan use mixes*

Table 1: Hypotheses related to Study 2

Hypothesis	Term	Explanation
$H_1^a$ : Borrowers who share one loan purpose are likely to share another	<code>cycle(4)</code>	This captures the tendency of borrowers and loan uses to be cyclical, where people cluster around shared uses of loans.
$H_2^a$ : Borrowers exhibit preferences for minimal debt distinctiveness, sharing one loan use	<code>b1degree(1)</code>	This captures the tendency of borrowers to prefer minimal distinctiveness in their loan uses, preferring to hold only one distinct type to minimise their perceived risk.
$H_3^a$ : Younger borrowers have different loan use mixes than older borrowers	<code>b1cov("age")</code>	This captures the tendency for younger borrowers to choose different kinds of loans based on how differently they consume credit.
$H_4^a$ : Gender differences lead to different loan use mixes	<code>b1nodematch("gender")</code>	This captures the idea that there are systemic differences in what kinds of loan uses are preferred the genders due to their perceived riskiness and creditworthiness.

Following this section, we outline the methodology, explaining our dataset and network construction. To facilitate this, descriptive statistics are analysed for both the dataset and resulting network. Subsequently, we elaborate on our choice of models, discussing how network models can be used to conduct the study. Then, modelling results are shown and explained to evaluate our hypotheses. Simultaneously, the robustness of the study is evaluated through the models' diagnostics. Finally, we arrive at our conclusions, summarising the study, its results, and implications. Supporting items such as the source code and specific data processing steps are shown in section [A](#).

## 2 Dataset

The study utilises publicly-available data from Bondora, a European P2P lending platform primarily operating in the Baltics and Spain. The dataset contains detailed information on defaulted and non-defaulted loans granted to users between February 2009 and July 2021. Specifically, contained is a range of numeric, binary, categorical, and time-series attributes across 85,087 unique users and 179,235 individual loans. Since the company's API is no longer accessible, we utilised a publicly available repository compiled by Manu Siddhartha (2021) on kaggle.com. The user made a series of API calls to collect the data. Following this, the dataset was processed according to Appendix [A.1](#).

## 2.1 Pre-Processing & Network Formation

Initially, the dataset was filtered loans between 2014 and 2016 since the data is incomplete during other periods. Since we are studying interactions across two distinct set of nodes, a bipartite network was formed with unique users as the first partition and the purpose of loans as the second. This allows us to evaluate how borrowers interact with loan purposes; a one-mode projection onto users is not meaningful for P2P networks as it produces no discernible structure. For model efficiency and performance, we randomly sample 500 individuals from the larger sample. Ultimately, we have 500 and 9 nodes in the first and second partitions, respectively.

## 2.2 Descriptive Statistics & Preliminary Analysis

Table 2 summarises the numeric variables between the reduced and larger samples. First, borrowers take out relatively small loans with fairly high interest rates, possibly reflecting that the average borrower is deemed fairly risky. The average loan is quite long, at approximately 46 months out of a maximum 60. Figure 9 shows that the largest credit category is HR<sup>1</sup>, which is the lowest possible rating. So, the typical borrower is relatively uncreditworthy. Moreover, the largest categories of loan use are “Other” and “Home Improvement”, suggesting these are largely personal. We observe minimal differences between the random sampling and the larger sample of individuals, per Table 2 and Figure 9, however, the study may be biased by the fact that the dataset contains only granted loans, rather than all bids for loans. This means our inferences about how borrowers behave is limited to successful bids. Further, the time period is relatively distant, meaning that recent structural changes to the P2P sector cannot be accounted for.

Table 2: Descriptive Statistics of Numeric Variables

DataFrame	Variable	Mean	Median	Std. Dev.	Min	Max
Initial Sample	Amount	2652.01	2125	2151.28	115	10630
Initial Sample	Interest	35.94	31	24.77	7.62	263.63
Initial Sample	Age	38.53	37	11.4	19	70
Initial Sample	LoanDuration	45.5	60	17.53	3	60
Reduced Sample	Amount	2637.88	2125	2110.65	170	10630
Reduced Sample	Interest	35.84	30	27.22	10.17	253.08
Reduced Sample	Age	38.73	37	11.33	21	69
Reduced Sample	LoanDuration	46.25	60	17.34	3	60

Table 3 shows a network density of 0.12, indicating that only 12% of all possible borrower–loan-type connections exist. This density suggests that borrowers typically engage with few loan types, reflecting specialization in borrowing behaviour. The centralization value of 0.47 further indicates an uneven distribution across loan types, where a few popular categories attract many borrowers while most remain peripheral. The mean distance of approximately three implies that borrowers are closely connected, typically separated by only three-to-four loan-use steps.

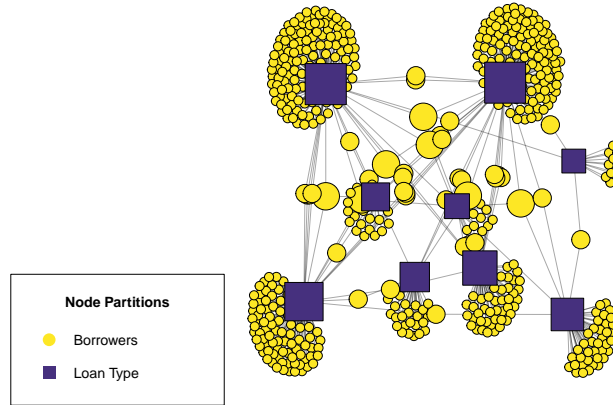
<sup>1</sup>Bondora uses their own proprietary credit rating system where the best rating is AA and worst is HR (Bondora, 2024). The company bases these ratings on their calculated expected probability of loss on the loan. For example, someone rated AA ranges from a expected loss of 0% to 2%, whereas for HR this is 25% - >25%.

Centrality distributions reinforce this structure Figure 2. Betweenness is highest among loan-type nodes, though several borrowers also bridge distinct loan clusters. Closeness centrality remains low but scattered, suggesting localized clustering constrained by the bipartite structure. Degree distributions show wide variation among loan types, consistent with a hub-like structure dominated by a few popular uses. However, the second partition has relatively small degree dispersion. Finally, the homogeneous eccentricity distribution indicates that no borrowers are highly isolated, reflecting a compact and moderately cohesive network overall.

Table 3: Basic Network Descriptive Statistics

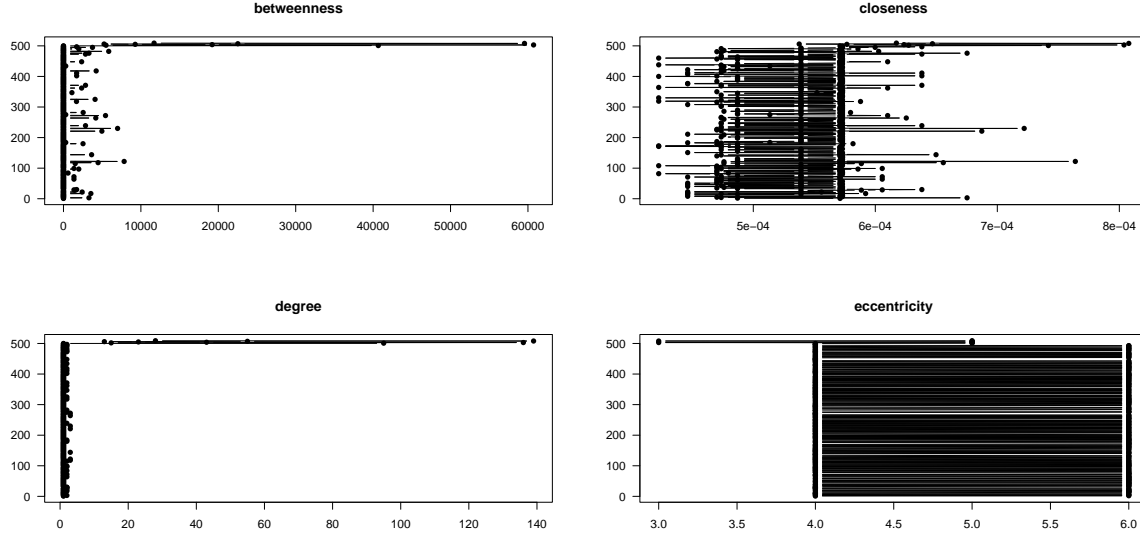
Statistic	Measure
Vertex Count	509.00
Edge Count	547.00
Density	0.12
Centralization	0.47
Mean Distance	3.66

Figure 1: Network Visualisation



Note: Each colour represents a different partition of the network. The partitions' nodes are scaled to their relative degrees. Users with many different loan types are larger than those with single loan uses. Similarly, the size of second partition nodes represents how many users belong to each.

Figure 2: Network Centrality Distribution Plots



### 3 Research Rationale

Our overall study analyses borrowers' behaviours, however, the studies are differentiated in their scope; the first study observes dyad-level relations among borrowers to understand how shared attributes affect borrowers' behaviours. This is inherently useful because behavioural finance aims to understand how individuals behave in *aggregate* (Ayal et al., 2018). Typically, the Quadratic Assignment Procedure is selected for dyad-level observations because basic linear regressions such as Ordinary Least Squares are fundamentally flawed and biased when using non-independent observations (Simpson, 2001).

Recognising this, studies such as Zuo (n.d.) assume that groups of borrowers have dependencies, and accordingly, use clustering methods such as K-Means or Fuzzy Clustering to understand behaviour. However, these methods are unsuitable when we require a distribution of estimates to make inferences from our results. Moreover, studying individuals as dyads rather than large clusters provides finer granularity in results. Therefore, QAP regression is the most appropriate choice for this study: it allows us to control for the non-independence among dyads, derive a statistical distribution of estimates through matrix permutation, and draw valid inferences about the relationship between borrower similarities and loan-use behaviour. This makes it both theoretically and methodologically well-suited to our research objectives and dataset.

The second study extends the analysis from a dyadic similarities to the structural formation of borrower-to-loan-use connections, which can identify the behavioural mechanisms that influence borrowers' specific loan use mixes. While QAP can account for exogenous influences, it cannot analyse higher-order network patterns such as clustering while accounting for simultaneous interactions between borrowers and loan uses. Using the Exponential Random Graph Model allows us to analyse how local behavioural tendencies aggregates to a global network structure.

Recent evidence highlights the need to focus on structural approaches in P2P lending. Liu, Baals,



Osterrieder, & Hadji-Misheva (2024) studies borrowers' centrality with the Bondora P2P dataset, modelling edges through borrower-to-borrower similarity metrics. Ultimately, they find that a borrower's position is highly significant for their probability of default, using centrality as an exogenous variable in a logistic regression. In contrast, ERGM allows us to move away from predictive association and instead account for endogenous factors, allowing us to understand the likelihood that borrowers' selection of loan uses are random. Consequently, we can better understand how behavioural phenomena such as familiarity and debt account aversion can translate into *systemic* patterns of credit use.

As we aim to understand how two distinct groups of nodes interact, a bipartite approach is highly suitable. Specifically, Stivala, Wang, & Lomi (2025) shows that in bipartite networks, terms such as `cycle(4)` can be effective in capturing clustering as in hypothesis 3, whereas it would fail in one-mode projected networks. Furthermore, the flexibility of ERGM also allows for modelling exogenous impacts such as heterophily in gender-based loan use mixes, as in hypothesis 6, and covariate age effects as in hypothesis 5.

## A Supplements

### A.1 Data Preprocessing Steps

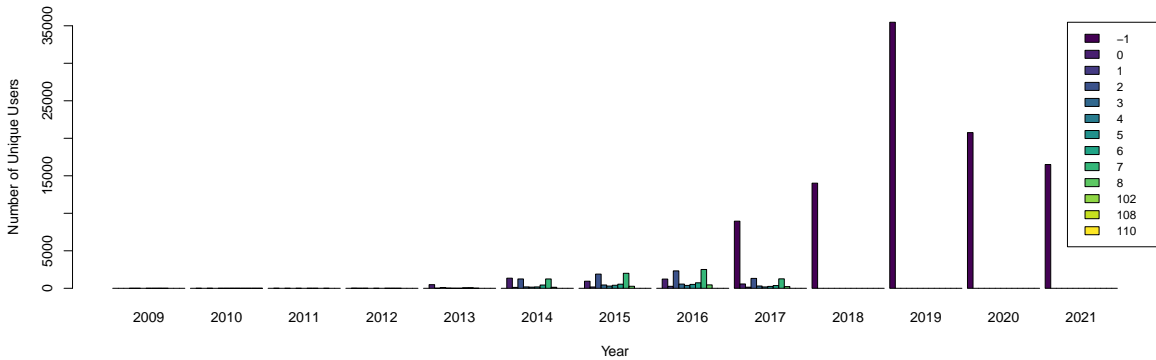
This section outlines the specific steps taken to process the data into the network object that was used in later analysis, outlining the rationale, where needed.

- Following the data import, only the following attributes were kept to make our processing steps more focused

```
keep_cols <- c("LoanId", "UserName", "Age", "Gender",  
              "Country", "Amount", "Interest", "LoanDuration",  
              "UseOfLoan", "Rating", "Restructured", "MonthlyPayment",  
              "OccupationArea", "BiddingStartedOn")
```

- We removed rows with any NA values for the selected attributes to ensure that we have a complete dataset
- Observing Figure 3, we see that between 2017-2021, there exists many -1 values for UseofLoan. This is because Bondora no longer provided this detail in later datasets. However, between 2014 and 2016 this data is available. Because of this, we restrict the time period to have the least noisy data. BiddingStartedOn was chosen to record the time, because the auction start date can be a good indication of activity around the loan.

Figure 3: Number of Unique Users per Loan Type per Year



- Then, we sampled 500 individuals out of the filtered dataset to ensure that our models run sufficiently well and converge. By choosing 500 individuals, we preserve importance structures within the network while adding sufficient statistical power.
- When adding labels to each loan use and occupation area, we keep any possible missing values because these do not indicate lower quality data. Instead, they can carry information about lenders distribute loans according to the quality of information presented by borrowers (Yao et al., 2019). Additionally, keeping missing labels is important to ensure that the edges are not superficially limited and that nodes can be isolates if they are indeed that way in reality.
- To create the network object for ERGM models:

- We first create an incidence matrix of size  $n \times m$ , borrower usernames by the loan use, where the values are the number of loans the user obtained for each loan use category. These values are then turned into binary because our modelling does not support the use of weighted edges.
- A bipartite network is created from this incidence matrix, where the first partition are unique users, and the second partition is the loan use associated among all loans of the user. Users can share more than more loan use type.
- We add age and gender attributes to each node of the first partition
- To create the relevant networks/matrices for QAP Regressions:
  - We first re-create the loan use incidence matrix, however, we keep the counts as the linear regression supports this.
    - \* We then convert the loan use incidence matrix into an adjacency matrix of shape  $n \times n$  through the transformation  $\mathbf{X} \cdot \mathbf{X}^T$ , which results in a matrix of shared loan uses weighted by the frequency of each user’s loan count within each loan use category.
    - \* We then set the diagonal values of  $\mathbf{X}$  to be zero to ensure that there are no loops within the network.
  - In accordance with this procedure, we create adjacency matrices for Credit Rating and Occupation Area which will be used as the main predictors in the QAP models.
    - \* Credit Rating is a weighted adjacency matrix based on the loan count belonging to each user for each credit rating category
    - \* Occupation Area is a binary adjacency matrix based on which occupation each user has reported to have held in the past when applying for each of their loans.
  - Additionally, we create a set of control variables to improve the validity and performance of the linear regression using Loan Amounts, Age, Gender, Loan Duration, and whether the loans have been restructured.
    - \* Loan Amounts is a weighted adjacency matrix based on five bins across the range of possible loan amounts
    - \* Age is an adjacency matrix based on differences in users’ ages
    - \* Gender is a binary adjacency matrix based on users’ shared gender.
    - \* Loan Duration is an adjacency matrix based on the differences in users’ loans average durations
    - \* Restructured is a binary adjacency matrix based on whether the users share the fact that they have defaulted on any loans in the past. The users can either have shared both defaults and non-defaults, or both.

A.2 Distributions of Variables across Samples

Figure 4: Distribution of Loan Amount across Samples

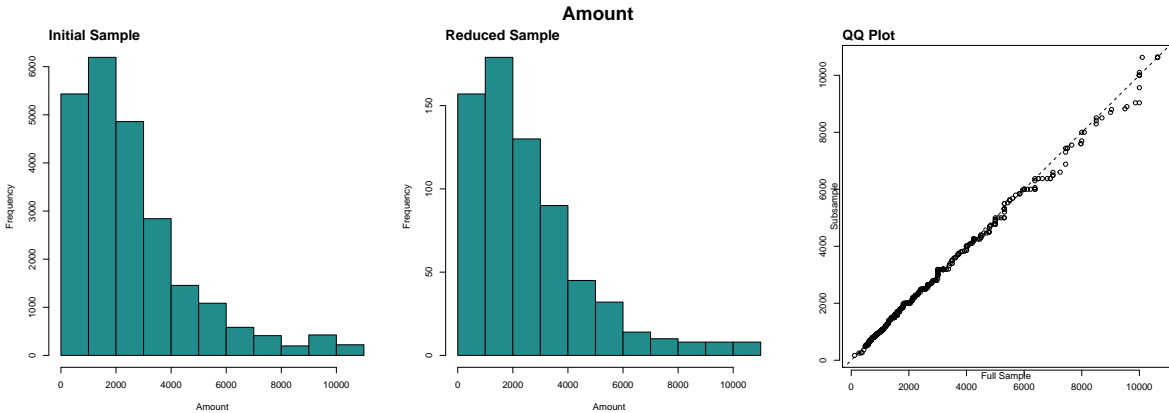


Figure 5: Distribution of Interest Rates across Samples

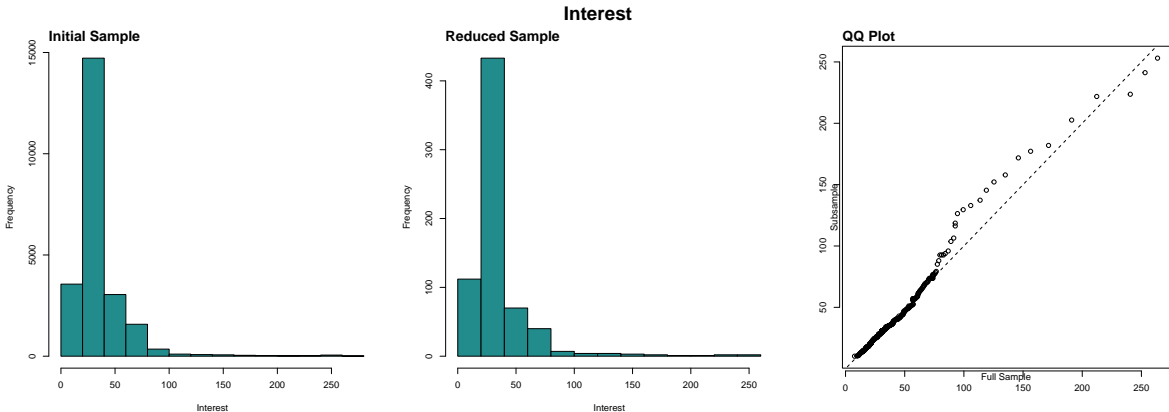


Figure 6: Distribution of Age across Samples

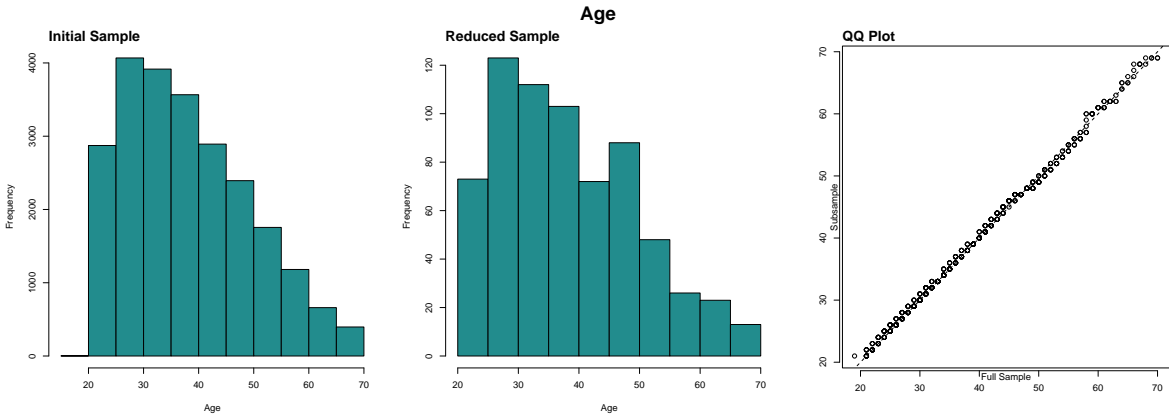


Figure 7: Distribution of Loan Duration across Samples

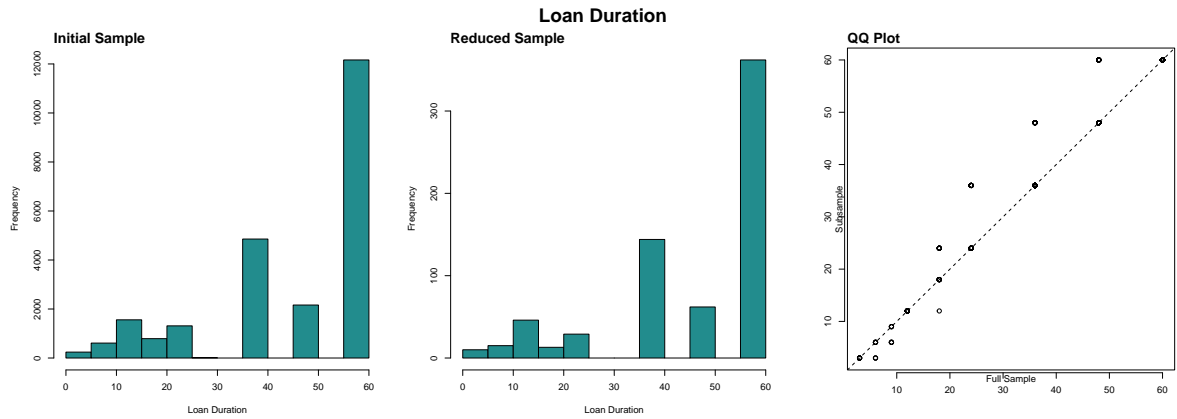


Figure 8: Distribution of Loan Purpose across Samples

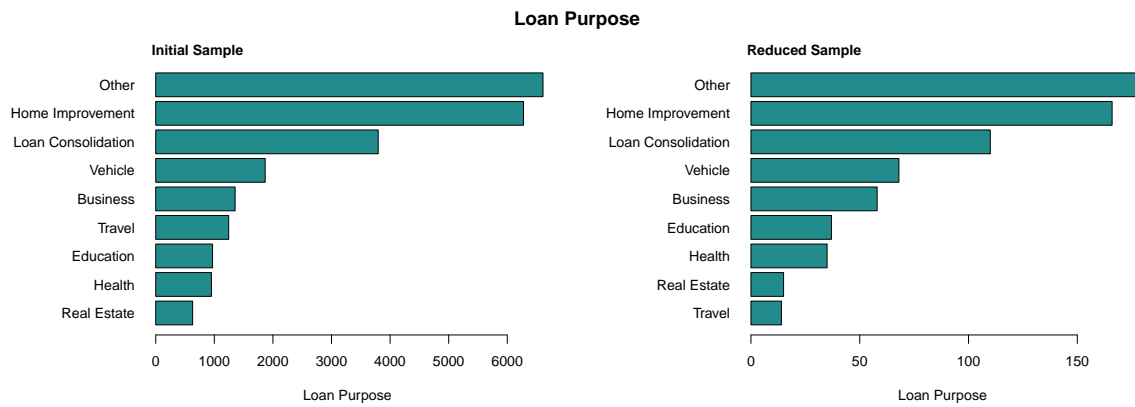


Figure 9: Distribution of Credit Ratings across Samples

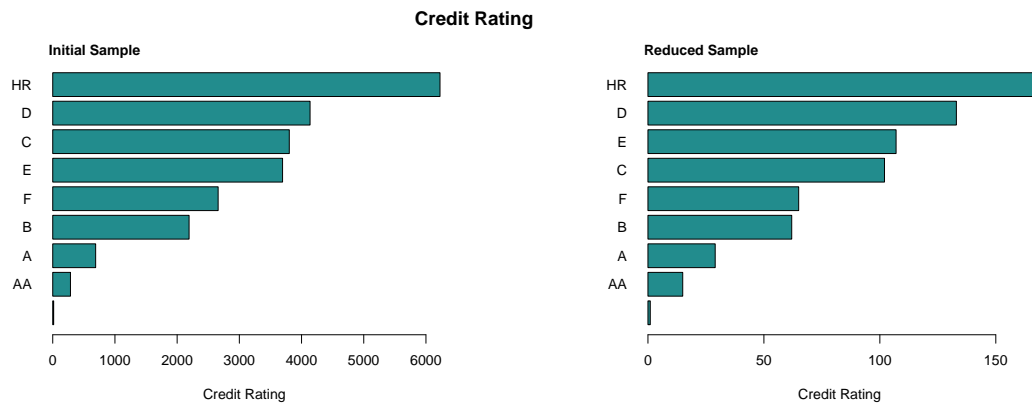
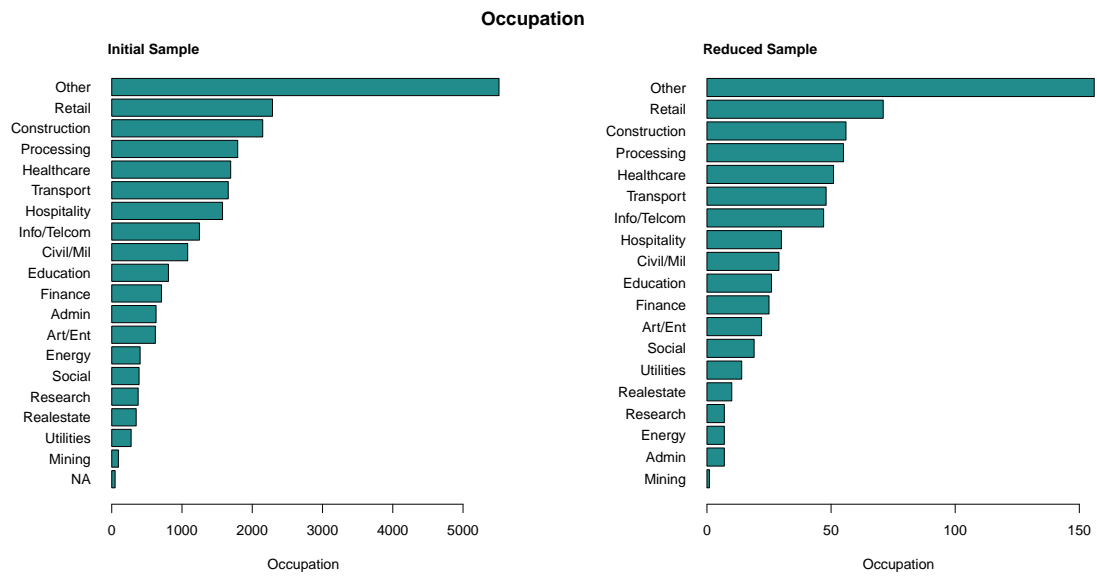


Figure 10: Distribution of Occupation Types across Samples



### A.3 Source Code - Data Preprocessing

```
1 # ----- #
2 install.packages("viridis") # For Colours
3 install.packages("here") # To locate files from RProj
4 set.seed(42)
5 # ----- #
6 # Import the Bondora P2P Dataset
7 obj_paths = "resources/objects/"
8 bondora_raw <- read.csv("dataset/LoanData_Bondora.csv")
9 raw_cols <- colnames(bondora_raw)
10 # ----- #
11 # Select Columns to Keep
12 keep_cols <- c("LoanId", "UserName", "Age", "Gender",
13               "Country", "Amount", "Interest", "LoanDuration",
14               "UseOfLoan", "Rating", "Restructured", "MonthlyPayment",
15               "OccupationArea", "BiddingStartedOn")
16
17 bondora <- bondora_raw[keep_cols]
18
19 # Change date format to the correct one
20 bondora$BiddingStartedOn <- as.POSIXct(bondora$BiddingStartedOn,
21                                       format = "%Y-%m-%d %H:%M:%S")
22 bondora$year <- as.numeric(format(bondora$BiddingStartedOn, "%Y"))
23
24 # Remove Rows with any NAs -> Complete Dataset Preferred
25 print(paste("NA Count |", sum(is.na(bondora)), "rows"))
26 bondora <- na.omit(bondora)
27 print(paste("NA Count |", sum(is.na(bondora)), "rows"))
28
29 # Observe the distribution of Loan Use over Time
30 unique_counts <- tapply(bondora$UserName,
31                         list(bondora$year, bondora$UseOfLoan),
32                         function(x) length(unique(x)))
33 unique_counts[is.na(unique_counts)] <- 0
34
35 barplot(t(unique_counts), # transpose so bars = loan types
36         beside = TRUE, # side-by-side bars
37         col = viridis::viridis(ncol(unique_counts)),
38         legend.text = colnames(unique_counts),
39         args.legend = list(x = "topright", cex = 0.8),
40         xlab = "Year",
41         ylab = "Number of Unique Users")
42 user_counts_plt <- recordPlot()
43 saveRDS(user_counts_plt,
44         here::here("resources", "objects", "preprocessing", "user_cnt_plt.Rds"))
45
46 # Restrict to Recent Time Period
47 bondora_test <- bondora[bondora$year > 2013 & bondora$year < 2017, ]
48
49 # See the distribution of LoanUses
50 see_counts <- function(var1, var2) {
51   counts <- table(var1, var2)
52   total_counts <- colSums(counts)
53   prop <- total_counts / sum(total_counts)
54
55   return(list(total_counts, round(prop, 2)))
56 }
57 before_counts <- see_counts(bondora_test$UserName, bondora_test$UseOfLoan)
58 barplot(before_counts[[1]])
59
60 # Number of unique users in the subsample
61 unique_users <- unique(bondora_test$UserName)
62 sample_users <- sample(unique_users, 500)
63 bondora_clean <- bondora_test[bondora_test$UserName %in% sample_users, ]
64
65 after_counts <- see_counts(bondora_clean$UserName, bondora_clean$UseOfLoan)
```

```

66 | barplot(after_counts[[1]])
67 |
68 | # DEPRECATED METHOD
69 | # Filter People with 2+ Loans
70 | #user_counts <- table(bondora_test$UserName)
71 | #multi_users <- names(user_counts[user_counts > 2])
72 | #bondora_clean <- bondora_test[bondora_test$UserName %in% multi_users, ]
73 |
74 | # DEPRECATED METHOD
75 | # Remove Users with only One Loan
76 | #user_counts <- table(bondora_clean$UserName)
77 | #multi_users <- names(user_counts[user_counts > 5])
78 | #bondora_clean <- bondora_clean[bondora_clean$UserName %in% multi_users, ]
79 | #bondora_test <- bondora_raw[bondora_raw$UserName %in% multi_users, ]
80 |
81 | # See if Ratings are Properly Encoded
82 | unique(bondora_clean$Rating)
83 |
84 | # Extract UseofLoan Types and Turn into Factor
85 | bondora_clean$UseOfLoan_factor <- as.factor(bondora_clean$UseOfLoan)
86 | unique(bondora_clean$UseOfLoan_factor)
87 |
88 | loan_use_labels <- c(
89 |   "-1" = "NA",
90 |   "0" = "Loan Consolidation",
91 |   "1" = "Real Estate",
92 |   "2" = "Home Improvement",
93 |   "3" = "Business",
94 |   "4" = "Education",
95 |   "5" = "Travel",
96 |   "6" = "Vehicle",
97 |   "7" = "Other",
98 |   "8" = "Health",
99 |   "110" = "Other Business",
100 |   "102" = "Undefined Business",
101 |   "108" = "Undefined Business"
102 | )
103 | # Change Labels for Cleaned Dataset
104 | bondora_clean$UseOfLoan_factor <- loan_use_labels[
105 |   as.character(bondora_clean$UseOfLoan)]
106 |
107 | # Change Labels for Uncleaned Dataset
108 | bondora_test$UseOfLoan_factor <- loan_use_labels[
109 |   as.character(bondora_test$UseOfLoan)]
110 |
111 | # Add labels to the OccupationArea Variable
112 | levels(as.factor(bondora_clean$OccupationArea)) # view codes
113 |
114 | occupation_labels <- c(
115 |   "-1" = "NA",
116 |   "1" = "Other",
117 |   "2" = "Mining",
118 |   "3" = "Processing",
119 |   "4" = "Energy",
120 |   "5" = "Utilities",
121 |   "6" = "Construction",
122 |   "7" = "Retail",
123 |   "8" = "Transport",
124 |   "9" = "Hospitality",
125 |   "10" = "Info/Telcom",
126 |   "11" = "Finance",
127 |   "12" = "Realestate",
128 |   "13" = "Research",
129 |   "14" = "Admin",
130 |   "15" = "Civil/Mil",
131 |   "16" = "Education",

```



```

132 "17" = "Healthcare",
133 "18" = "Social",
134 "19" = "Art/Ent",
135 "20" = "Agriculture",
136 "21" = "Forestry/Fish"
137 )
138 # store original
139 bondora_clean$occupation_code <- bondora_clean$OccupationArea
140 bondora_clean$occupation_label <- occupation_labels[
141   as.character(bondora_clean$OccupationArea)]
142
143 bondora_test$occupation_code <- bondora_test$OccupationArea
144 bondora_test$occupation_label <- occupation_labels[
145   as.character(bondora_test$OccupationArea)]
146 # ----- #
147 # Observe Descriptive Statistics
148
149 cols <- viridis::viridis(30)
150
151 # Make function to save plots
152 save_plot <- function(plt_nam) {
153   plt <- recordPlot()
154   saveRDS(plt, here::here("resources", "objects", "preprocessing",
155     paste0(plt_nam, ".Rds")))
156 }
157
158 # Make function to consistently plot comparisons
159 plot_desc_hists <- function(df1, df2, col_name, type) {
160
161   par(mfrow=c(1,3))
162
163   hist(df1[[col_name]], xlab=type, col=cols[15], main="", breaks=10)
164   mtext("Initial Sample", side=3, adj=0, line=0.25, cex=1, font=2)
165
166   hist(df2[[col_name]], xlab=type, col=cols[15], main="", breaks=10)
167   mtext("Reduced Sample", side=3, adj=0, line=0.25, cex=1, font=2)
168
169   qqplot(df1[[col_name]], df2[[col_name]], main="", cex=1,
170     xlab="Full Sample", ylab="Subsample", line=0.25)
171   abline(0, 1, lty=2)
172
173   mtext("QQ Plot", side=3, adj=0, line=0.25, cex=1, font=2)
174
175   mtext(type, outer = TRUE, line = -2, side=3, cex = 1.3, font = 2)
176
177   # Reset plot window
178   par(mfrow=c(1,1), mar=c(5,4,4,2)+0.1)
179 }
180
181 plot_desc_bar <- function(df1, df2, col_name, type) {
182
183   par(mfrow=c(1,2), mar=c(5,10,4,2))
184
185   barplot(sort(table(df1[[col_name]]), decreasing = F),
186     xlab=type, col=cols[15], horiz=TRUE, las=1)
187   mtext("Initial Sample", side=3, adj=0, line=0.25, cex=1, font=2)
188
189   barplot(sort(table(df2[[col_name]]), decreasing = F),
190     xlab=type, col=cols[15], horiz=TRUE, las=1)
191   mtext("Reduced Sample", side=3, adj=0, line=0.25, cex=1, font=2)
192
193   mtext(type, outer = TRUE, line = -2, side=3, cex = 1.3, font = 2)
194
195   # Reset plot window
196   par(mfrow=c(1,1), mar=c(5,4,4,2)+0.1)
197 }

```

```

198
199 plot_desc_hists(bondora_test, bondora_clean, "Amount", "Amount")
200 save_plot("hist_amt")
201
202 plot_desc_hists(bondora_test, bondora_clean, "Interest", "Interest")
203 save_plot("hist_int")
204
205 plot_desc_hists(bondora_test, bondora_clean, "LoanDuration", "Loan Duration")
206 save_plot("hist_loandur")
207
208 plot_desc_hists(bondora_test, bondora_clean, "MonthlyPayment", "Monthly Payment")
209 save_plot("hist_monpmt")
210
211 plot_desc_hists(bondora_test, bondora_clean, "Age", "Age")
212 save_plot("hist_age")
213
214 plot_desc_bar(bondora_test, bondora_clean, "UseOfLoan_factor", "Loan Purpose")
215 save_plot("bar_loanuse")
216
217 plot_desc_bar(bondora_test, bondora_clean, "Rating", "Credit Rating")
218 save_plot("bar_rating")
219
220 plot_desc_bar(bondora_test, bondora_clean, "occupation_label", "Occupation")
221 save_plot("bar_occupation")
222
223 # Get Tabular Summary Statistics
224 tab_comps <- function(df1, df2, cols) {
225   stats <- c("Mean"=mean, "Median"=median, "Std. Dev."=sd, "Min"=min, "Max"=max)
226
227   get_stats <- function(d) {
228     t(sapply(d[cols], function(x)
229       sapply(stats, function(f) round(f(x, na.rm=TRUE), 2))
230     ))
231   }
232
233   df1_stats <- get_stats(df1)
234   df2_stats <- get_stats(df2)
235
236   out <- rbind(
237     cbind(DataFrame = "Initial Sample",
238       Variable = rownames(df1_stats), df1_stats),
239     cbind(DataFrame = "Reduced Sample",
240       Variable = rownames(df2_stats), df2_stats)
241   )
242   rownames(out) <- NULL
243   as.data.frame(out)
244 }
245
246 tab_results <- tab_comps(bondora_test, bondora_clean,
247   c("Amount", "Interest", "Age", "LoanDuration"))
248
249 knitr::kable(tab_results)
250
251 # Save table for use in the Report
252 saveRDS(tab_results,
253   file=paste0(obj_paths, "preprocessing/", "summary_table.Rds"))
254 # ----- #
255 # Convert Dataset into Incidence Matrix to form Network Object (for ERGM)
256 bondora_slim <- bondora_clean
257
258 # Create the Incidence Matrix for Use of Loan
259 bondora_matrix <- table(
260   bondora_slim$UserName, bondora_slim$UseOfLoan)
261 bondora_matrix[bondora_matrix > 0] <- 1 # Given that ergm.counts fails with GOF
262
263 # Create network object with counts as Edge attribute

```

```

264 bondora_net <- network::network(
265   bondora_matrix, directed=FALSE, bipartite=nrow(bondora_matrix),
266   ignore.eval = FALSE, names.eval="frequency", loops=FALSE)
267
268 # Set the bipartite Attribute - UNNECESSARY GIVEN bipartite=length(n)
269 len <- dim(bondora_matrix)[1]
270 len_b2 <- dim(bondora_matrix)[2]
271 b_indicator <- c(rep(1,len),rep(2,len_b2))
272
273 # Extract Partition 2 Labels
274 loan_use <- levels(bondora_clean$UseOfLoan_factor)
275
276 # Create Loan Type Attribute for Partition 2
277 b2_loantype <- rep(NA, len)
278 b2_loantype <- c(b2_loantype, loan_use)
279
280 if (length(b2_loantype) == network::network.size(bondora_net)) {
281   network::set.vertex.attribute(
282     bondora_net, "b2_loantype", value = b2_loantype)
283 }
284
285 # Add Age Vertex Attribute to B1
286 age <- bondora_clean$Age[match(
287   rownames(bondora_matrix), bondora_clean$UserName)]
288 b1_age <- c(age, rep(NA, len_b2))
289 network::set.vertex.attribute(
290   bondora_net, "b1_age", value=b1_age
291 )
292
293 # Add Gender Vertex Attribute to B1
294 gender <- bondora_clean$Gender[match(
295   rownames(bondora_matrix), bondora_clean$UserName)]
296 unique(gender) # Check to see if encoded properly
297 gender <- ifelse(gender == 0, "male", "female")
298 b1_gender <- c(gender, rep(NA, len_b2))
299 network::set.vertex.attribute(
300   bondora_net, "b1_gender", value = b1_gender
301 )
302
303 # Save the network and data frame object
304 saveRDS(bondora_slim, file=paste0(obj_paths, "preprocessing/", "bondora_df.Rds"))
305 saveRDS(bondora_net, file=paste0(obj_paths, "preprocessing/", "bondora_net.Rds"))
306 # ----- #
307 # First, create Incidence Matrix again but with Frequency Counts
308 loan_use_matrix <- table(
309   bondora_slim$UserName, bondora_slim$UseOfLoan)
310
311 # Get the Adjacency Matrix for Loan Use Similarity (Dependent QAP Variable)
312 adj_mat_loan_use <- loan_use_matrix %>% t(loan_use_matrix)
313 # Remove self weights to remove any loops
314 diag(adj_mat_loan_use) <- 0
315
316 # Get the Adjacency Matrix for Credit Rating Similarity (Predictor in QAP)
317 incidence_rating <- table(bondora_slim$UserName, bondora_slim$Rating)
318 adj_mat_rating <- incidence_rating %>% t(incidence_rating)
319 diag(adj_mat_rating) <- 0
320
321 # Get Adjacency Matrix for Occupation Similarity (Predictor in QAP, Binary)
322 incidence_occupation <- table(bondora_slim$UserName,
323   bondora_slim$occupation_label)
324 adj_mat_occupation <- incidence_occupation %>% t(incidence_occupation)
325 adj_mat_occupation[adj_mat_occupation > 0] <- 1
326 diag(adj_mat_occupation) <- 0
327
328 # Get the Adjacency Matrix for Loan Amount (Control in QAP) - BINARY
329 # First, bin the Loan Amounts

```

```

330 bondora_slim$Amount_bins <- cut(
331   bondora_slim$Amount, breaks=c(0,2000,4000,6000,8000,10000),
332   labels = c(1:5)
333 )
334 incidence_amount_bins <- table(bondora_slim$UserName, bondora_slim$Amount_bins)
335 adj_mat_amount_bins <- incidence_amount_bins %>% t(incidence_amount_bins)
336 diag(adj_mat_amount_bins) <- 0
337
338 # Continous Absolute Difference Approach
339 avg_amount <- tapply(bondora_slim$Amount, bondora_slim$UserName, mean)
340 adj_mat_amount_diff <- outer(avg_amount, avg_amount,
341   FUN = function(x,y) abs(x - y))
342 diag(adj_mat_amount_diff) <- 0
343
344 # Get Matrix for Differences in Age
345 incidence_age <- table(bondora_slim$UserName, bondora_slim$Age)
346 borrower_ages <- as.numeric(colnames(incidence_age)[max.col(incidence_age)])
347 names(borrower_ages) <- rownames(incidence_age)
348 adj_mat_age <- outer(borrower_ages, borrower_ages,
349   FUN = function(x, y) abs(x - y))
350 rownames(adj_mat_age) <- colnames(adj_mat_age) <- names(borrower_ages)
351
352 # Get Adjacency Matrix for (same) Gender
353 incidence_gender <- table(bondora_slim$UserName, bondora_slim$Gender)
354 adj_mat_gender <- incidence_gender %>% t(incidence_gender)
355 # Make the matrix binary for homophily
356 adj_mat_gender <- ifelse(adj_mat_gender > 0, 1, 0)
357 diag(adj_mat_gender) <- 0
358
359 # Get Adjacency Matrix for Differences in Average Loan Duration
360 borrower_loandur <- sapply(tapply(bondora_slim$LoanDuration,
361   bondora_slim$UserName, unique),
362   mean)
363 adj_mat_loandur_diff <- outer(borrower_loandur, borrower_loandur,
364   FUN=function(x,y) abs(x-y))
365 diag(adj_mat_loandur_diff) <- 0
366
367 # Get Adjacency Matrix for Homophily in Restructure of Loans
368 incidence_restructure <- table(bondora_slim$UserName, bondora_slim$Restructured)
369 adj_mat_rest <- incidence_restructure %>% t(incidence_restructure)
370 diag(adj_mat_rest) <- 0
371
372 # Save the objects for the QAP Regression in different Script
373 qap_paths = paste0(obj_paths, "/qap/")
374 saveRDS(b_indicator, file=paste0(
375   "resources/objects/preprocessing/indicator.Rds"))
376
377 saveRDS(adj_mat_loan_use, file=paste0(qap_paths, "adj_mat_loanuse.Rds"))
378 saveRDS(adj_mat_occupation, file=paste0(qap_paths, "adj_mat_occup.Rds"))
379 saveRDS(adj_mat_rating, file=paste0(qap_paths, "adj_mat_rating.Rds"))
380 saveRDS(adj_mat_amount_diff, file=paste0(qap_paths, "adj_mat_amtdiffs.Rds"))
381 saveRDS(adj_mat_age, file=paste0(qap_paths, "adj_mat_agediffs.Rds"))
382 saveRDS(adj_mat_gender, file=paste0(qap_paths, "adj_mat_gender.Rds"))
383 saveRDS(adj_mat_loandur_diff, file=paste0(qap_paths, "adj_mat_loandurdiffs.Rds"))
384 saveRDS(adj_mat_rest, file=paste0(qap_paths, "adj_mat_rest.Rds"))
385 # ----- #

```

scripts/bondora\_preprocessing.R

## A.4 Source Code - QAP Linear Regression

```

1 # ----- #
2 # NOTE: NOT THE LATEST QAP ANALYSIS!
3
4 # If not already Installed
5 install.packages("viridis") # For Colours
6 install.packages("here") # To locate files from RProj
7
8 # Set colour palette
9 cols <- viridis::viridis(30)
10 # ----- #
11 # Load Relevant Files
12 qap_path="resources/objects/qap/"
13
14 loan_use_mat <- readRDS(paste0(qap_path, "adj_mat_loanuse.RDS"))
15 rating_mat <- readRDS(paste0(qap_path, "adj_mat_rating.RDS"))
16 amt_diffs_mat <- readRDS(paste0(qap_path, "adj_mat_amtdiffs.RDS"))
17 age_diffs_mat <- readRDS(paste0(qap_path, "adj_mat_agediffs.RDS"))
18 gender_mat <- readRDS(paste0(qap_path, "adj_mat_gender.RDS"))
19 loandur_diffs_mat <- readRDS(paste0(qap_path, "adj_mat_loandurdiffs.RDS"))
20 rest_mat <- readRDS(paste0(qap_path, "adj_mat_rest.RDS"))
21 occup_mat <- readRDS(paste0(qap_path, "adj_mat_occup.Rds"))
22 # ----- #
23 # Create function to determine significance from t-value statistic
24 t_to_stars <- function(t) {
25   stars <- rep("", length(t))
26   stars[abs(t) >= 1.96] <- "*"
27   stars[abs(t) >= 2.576] <- "**"
28   stars[abs(t) >= 3.291] <- "***"
29
30   return(stars)
31 }
32
33 # Make function to save plots
34 save_qap_plot <- function(plt_nam) {
35   plt <- recordPlot()
36   saveRDS(plt, here::here("resources", "objects", "qap",
37     paste0(plt_nam, ".Rds")))
38 }
39
40 var_names <- c("Intercept", "Rating", "Occupation",
41   "Loan Amount", "Age", "Gender", "Loan Duration", "Restructured")
42 pred_vars <- list(rating_mat, occup_mat, amt_diffs_mat, age_diffs_mat,
43   gender_mat, loandur_diffs_mat, rest_mat)
44 # ----- #
45 # Basic QAP Linear Regression
46 qap_m1 <- sna::netlm(y = loan_use_mat,
47   x = list(rating_mat, amt_diffs_mat, age_diffs_mat,
48     gender_mat, loandur_diffs_mat, rest_mat),
49   nullhyp = "qapspp", reps = 2500)
50 qap_m1$names <- var_names
51 summary(qap_m1)
52
53 results_m1 <- qap_m1$coefficients
54 names(results_m1) <- var_names
55 results_sig_m1 <- paste(round(results_m1,3), t_to_stars(qap_m1$tstat))
56
57 # Plot Residuals
58 hist(qap_m1$residuals, main="QAP Residuals", col = cols[15])
59
60 # Save Model
61 saveRDS(qap_m1, file = paste0(qap_path, "qap_m1.RDS"))
62 # ----- #
63 # Standardised QAP Linear Regression
64 scaled_dep <- scale(loan_use_mat)
65 scaled_pred <- lapply(pred_vars, scale)

```

```

66 |
67 | qap_m2 <- sna::netlm(y = scaled_dep,
68 |                    x = scaled_pred,
69 |                    nullhyp = "qapspp", reps = 2500)
70 | qap_m2$names <- var_names
71 | summary(qap_m2)
72 |
73 | # Plot the result
74 | results_m2 <- qap_m2$coefficients
75 | names(results_m2) <- var_names
76 | results_sig_m2 <- paste(round(results_m2,3), t_to_stars(qap_m2$tstat))
77 |
78 | # Plot Residuals
79 | hist(qap_m2$residuals, main="QAP Residuals", col = cols[15])
80 |
81 | # Save Model
82 | saveRDS(qap_m2, file = paste0(qap_path, "qap_m2.RDS"))
83 | # ----- #
84 | # Plot the result
85 | par(mfrow=c(1,2))
86 |
87 | qap_plot_m1 <- barplot(results_m1, col = cols[15], border = cols[10],
88 |                      ylim = c(min(results_m1) + min(results_m1)*0.15,
89 |                               max(results_m1) + max(results_m1)*0.15),
90 |                      main="QAP Model Results (Unstandardised)")
91 | text(x = qap_plot_m1,
92 |      y = results_m1 + sign(results_m1)*(0.075*diff(range(results_m1))),
93 |      labels = results_sig_m1, font = 2)
94 |
95 |
96 | qap_plot_m2 <- barplot(results_m2, col = cols[15], border = cols[10],
97 |                      ylim = c(min(results_m2) + min(results_m2)*0.15,
98 |                               max(results_m2) + max(results_m2)*0.15),
99 |                      main="QAP Model Results (Standardised)")
100 | text(x = qap_plot_m2,
101 |      y = results_m2 + sign(results_m2)*(0.075*diff(range(results_m2))),
102 |      labels = results_sig_m2, font = 2)
103 |
104 | save_qap_plot("unstd_std_plot")
105 |
106 | par(mfrow=c(1,1))

```

scripts/qap\_network\_analysis.R

## A.5 Source Code - ERGM Network Analysis

```
1 # ----- #
2 # NOTE: NOT THE LATEST ERGM ANALYSIS! FOR REFERENCE ONLY
3
4 # If not already Installed
5 install.packages("viridis")      # For Colours
6 install.packages("Rglpk")        # Additional solver for ERGMs
7 install.packages("here")         # To locate files from RProj
8
9 # Import the Network and Other Object
10 ergm_path <- "resources/objects/ergm/"
11 bondora_net <- readRDS(here::here(
12   "resources", "objects", "preprocessing", "bondora_net.Rds"))
13 b_indicator <- readRDS(here::here(
14   "resources", "objects", "preprocessing", "indicator.Rds"))
15
16 # Set colour palette
17 cols <- viridis::viridis(30)
18
19 # Determine acceptable core count
20 n_cores <- parallel::detectCores() - 3 # Leave some out for other processes
21 print(paste("You have", n_cores, "usable cores"))
22
23 # Repeatability
24 set.seed(42)
25
26 # Save plots
27 save_ergm_plot <- function(plt_nam) {
28   plt <- recordPlot()
29   saveRDS(plt, here::here("resources", "objects", "ergm",
30     paste0(plt_nam, ".Rds")))
31 }
32 # ----- #
33 # Copy network for plotting
34 bondora_plot <- bondora_net
35
36 # Get node type for plotting
37 type_indicator <- ifelse(b_indicator == 2, TRUE, FALSE)
38 shape <- ifelse(type_indicator, "square", "circle")
39 network::set.vertex.attribute(bondora_plot, "shape", shape)
40
41 # Get Category Count for Vertex Size
42 counts <- sna::degree(bondora_plot)
43 counts_att <- ifelse(type_indicator, log(counts)*4, counts*2.5)
44 network::set.vertex.attribute(bondora_plot, "size", counts_att)
45
46 # Colours for the Node Types
47 plot_cols <- ifelse(type_indicator, cols[5], cols[30])
48 network::set.vertex.attribute(bondora_plot, "color", plot_cols)
49
50 # Legend Plotting
51 type_legend <- ifelse(type_indicator, "Borrowers", "Loan Type")
52 type_legend <- as.factor(type_legend)
53
54 # Plot the Network
55 plot(snafun::to_igraph(bondora_plot),
56   #main = "Bipartite User-LoanUse",
57   edge.arrow.size = 0.3,
58   edge.color = rgb(0,0,0, alpha = 0.35),
59   vertex.frame.color = "black",
60   vertex.label = NA,
61   vertex.frame.size = 3,
62   edge.curved = FALSE,
63   layout=igraph::layout.fruchterman.reingold)
64 legend("bottomleft",
65   legend = levels(type_legend),
```

```

66     inset = c(0.15, 0.01),
67     col = c(cols[30], cols[5]),
68     pch = c(16, 15),
69     title = "Node Partitions",
70     title.font = 2,
71     cex = 1,                # Increase the text size
72     pt.cex = 2,            # Increase the point symbol size
73     box.lwd = 1,           # Thin box border
74     box.col = "black",     # Box color
75     bty = "o"              # Use a box around legend
76 )
77 save_ergm_plot("network_plot")
78
79 # Summary Statistics and Save Them
80 density <- snafun::g_density(bondora_net)[1]
81 centralization <- snafun::g_centralize(bondora_net)[1]
82 vertices <- snafun::count_vertices(bondora_net)[1]
83 edges <- snafun::count_edges(bondora_net)[1]
84 dist <- snafun::g_mean_distance(bondora_net)[1]
85
86 net_names <- c("Vertex Count", "Edge Count", "Density", "Centralization",
87               "Mean Distance")
88 net_stats <- c(vertices, edges, density, centralization, dist)
89 net_stats <- sapply(net_stats, function(x) round(as.numeric(x), 2))
90
91 net_summary <- data.frame(Statistic = net_names,
92                           "Measure" = net_stats)
93 knitr::kable(net_summary)
94 saveRDS(net_summary, here::here("resources", "objects", "ergm", "net_summary.Rds"))
95
96 # Plot Network Summary Statistics
97 snafun::plot_centralities(bondora_net)
98 save_ergm_plot("network_plots")
99 # ----- #
100 # Make function to calculate probabilities from log odds
101 lodds_to_prob <- function(l_odd) {
102   return(exp(l_odd) / (1 + exp(l_odd)))
103 }
104 # Make function to save ERGM object
105 save_ergm <- function(object, id) {
106   saveRDS(object, file=here::here(
107     "resources", "objects", "ergm", id, ".Rds"))
108 }
109 # Make function to conduct ERGMs automatically
110 auto_ergm <- function(model, mcmc, name) {
111
112   # Conducts the GOF Diagnostics and then saves the model,
113   # mcmc diagnostics and gof object in a list.
114   # This list can be imported as an .RDS object into the R environment
115
116   # Diagnostics
117   if (mcmc) {
118     ergm::mcmc.diagnostics(model)
119   }
120
121   # The GOF must be adjusted otherwise it takes too long
122   # We do not limit the GOF by changing its range of parameters
123   gof <- ergm::gof(model,
124                     control = ergm::control.gof.ergm(
125                       nsim = 200,
126                       MCMC.burnin = 5000,
127                       MCMC.interval = 1000,
128                       parallel = n_cores,
129                       parallel.type = "PSOCK"
130                     ))
131

```



```

132 # Return List to view each item separately
133 result <- list(model, gof)
134 names(result) <- c("model", "gof")
135 save_ergm(result, paste0(name, "_panel"))
136
137 return(result)
138 }
139 # ----- #
140 # Find max degree
141 (max_deg <- max(summary(bondora_net ~ b2factor("b2_loantype"))))
142 # ----- #
143 # Base Model + GOF
144 formula_base_model <- bondora_net ~ edges
145 base_ergm <- ergm::ergm(formula_base_model)
146 base_ergm_panel <- auto_ergm(base_ergm, mcmc = FALSE, name = "ergm_base")
147 snafun::stat_plot_gof(base_ergm_panel$gof)
148 models = list(base_ergm)
149 texreg::screenreg(models)
150 # ----- #
151 # Base Model + Edge Counts + GOF
152 #base_model_counts <- ergm::ergm(bondora_net ~ edges, response="frequency",
153 #                                reference = ~ Poisson)
154 #basemodel_counts_gof <- ergm::gof(base_model_counts)
155 #snafun::stat_plot_gof(basemodel_counts_gof)
156 #
157 #texreg::screenreg(list(base_model, base_model_counts))
158 # ----- #
159 # Iteration 1 + MCMC Diagnostics + GOF
160 model_1_params <- bondora_net ~ edges +
161 # b1 decay can be very low since 9 b2
162 gwb1degree(decay=0.15, fixed=TRUE)
163
164 model_1 <- ergm::ergm(
165   model_1_params,
166
167   # Max b2 degree is 72, so this constraint is reasonable
168   # and helps convergence significantly.
169   # Technically in the Bondora population this can be
170   # far higher but we are studying a subsample.
171   #constraints = ~ bd(minout = 0, maxout = 80),
172
173   control = ergm::control.ergm(
174     # Greater burn-in for cleaner result
175     MCMC.burnin = 20000,
176     # Greater sample size for greater stability
177     MCMC.samplesize = 100000,
178     seed = 42,
179     MCMC.interval = 1000,
180     # Only needed for convergence pvals to improve
181     MCMLE.maxit = 45,
182     # Smaller steps for stability
183     MCMLE.steplength = 0.25,
184     parallel = n_cores,
185     parallel.type = "PSOCK"
186   )
187 )
188
189 model_1_panel <- auto_ergm(model=model_1, mcmc=TRUE, name="ergm_m1")
190 model_1_panel$gof
191 snafun::stat_plot_gof(model_1_panel$gof)
192 texreg::screenreg(list(base_ergm, model_1))
193 # ----- #
194 # Iteration 2 + MCMC Diagnostics + GOF
195 model_2_params <- bondora_net ~ edges +
196 # low decay important because there is high clustering around low degrees
197 gwb1degree(decay=0.15, fixed=TRUE) +

```

```

198 # decay should be higher due to wider variation in degree but
199 # too high of degree makes the traces concentrated around the tails.
200 gwbldsp(decay=0.5, fixed=TRUE)
201
202 model_2 <- ergm::ergm(
203   model_2_params,
204
205   # Max b2 degree is 72, so this constraint is reasonable
206   # and helps convergence significantly.
207   # Technically in the Bondora population this can be
208   # far higher but we are studying a subsample.
209   constraints = ~ bd(minout = 0, maxout = 80),
210
211   control = ergm::control.ergm(
212     # Greater burn-in for cleaner result
213     MCMC.burnin = 20000,
214     # Greater sample size for greater stability
215     MCMC.samplesize = 100000,
216     seed = 42,
217     MCMC.interval = 1000,
218     # Only needed for convergence pvals to improve
219     MCMLE.maxit = 45,
220     # Smaller steps for stability
221     MCMLE.steplength = 0.25,
222     parallel = n_cores,
223     parallel.type = "PSOCK"
224   )
225 )
226
227 model_2_panel <- auto_ergm(model=model_2, mcmc=TRUE, name="ergm_m2")
228 snafun::stat_plot_gof(model_2_panel$gof)
229 model_2_panel$gof
230 models <- list(base_ergm, model_1, model_2)
231 texreg::screenreg(models)
232 # ----- #
233 # Iteration 3 + MCMC Diagnostics + GOF
234 model_3_params <- bondora_net ~ edges +
235   # low decay important because there is high clustering around low degrees
236   gwbldegree(decay=0.15, fixed=TRUE) +
237   # decay should be higher due to wider variation in degree but
238   # too high of degree makes the traces concentrated around the tails.
239   gwbldsp(decay=0.5, fixed=TRUE) +
240   # See differences across genders (implicitly, since blnodemix unavailable)
241   blnodematch("bl_gender", diff=FALSE)
242
243 model_3 <- ergm::ergm(
244   model_3_params,
245
246   # Max b2 degree is 72, so this constraint is reasonable
247   # and helps convergence significantly.
248   # Technically in the Bondora population this can be
249   # far higher but we are studying a subsample.
250   constraints = ~ bd(minout = 0, maxout = 80),
251
252   control = ergm::control.ergm(
253     # Greater burn-in for cleaner result
254     MCMC.burnin = 20000,
255     # Greater sample size for greater stability
256     MCMC.samplesize = 100000,
257     seed = 42,
258     MCMC.interval = 1000,
259     # Only needed for convergence pvals to improve
260     MCMLE.maxit = 45,
261     # Smaller steps for stability
262     MCMLE.steplength = 0.25,
263     parallel = n_cores,

```

```

264     parallel.type = "PSOCK"
265   )
266 )
267
268 model_3_panel <- auto_ergm(model=model_3, mcmc=TRUE, name="ergm_m3")
269 snafun::stat_plot_gof(model_3_panel$gof)
270 model_3_panel$gof
271 models <- list(base_ergm, model_1, model_2, model_3)
272 texreg::screenreg(models)
273 # ----- #
274 # Iteration 4 + MCMC Diagnostics + GOF
275 model_4_params <- bondora_net ~ edges +
276   # low decay important because there is high clustering around low degrees
277   gwblddegree(decay=0.15, fixed=TRUE) +
278   # decay should be higher due to wider variation in degree but
279   # too high of degree makes the traces concentrated around the tails.
280   gwbldsp(decay=0.5, fixed=TRUE) +
281   # See differences across genders (implicitly, since blnodemix unavailable)
282   blnodematch("bl_gender", diff=TRUE) +
283   # See if higher ages make a difference
284   blcov("bl_age")
285
286 model_4 <- ergm::ergm(
287   model_4_params,
288
289   # Max b2 degree is 72, so this constraint is reasonable
290   # and helps convergence significantly.
291   # Technically in the Bondora population this can be
292   # far higher but we are studying a subsample.
293   constraints = ~ bd(minout = 0, maxout = 80),
294
295   control = ergm::control.ergm(
296     # Greater burn-in for cleaner result
297     MCMC.burnin = 20000,
298     # Greater sample size for greater stability
299     MCMC.samplesize = 100000,
300     seed = 42,
301     MCMC.interval = 1000,
302     # Only needed for convergence pvals to improve
303     MCMLE.maxit = 45,
304     # Smaller steps for stability
305     MCMLE.steplength = 0.25,
306     parallel = n_cores,
307     parallel.type = "PSOCK"
308   )
309 )
310
311 model_4_panel <- auto_ergm(model=model_4, mcmc=TRUE, name="ergm_m4")
312 model_4_panel$gof
313 snafun::stat_plot_gof(model_4_panel$gof)
314 models <- list(base_ergm, model_1, model_2, model_3, model_4)
315 texreg::screenreg(models)
316 # ----- #

```

scripts/ergm\_network\_analysis.R

## References

- Agrawal, V. (2012). Managing the diversified team: Challenges and strategies for improving performance. *Team Performance Management: An International Journal*, 18(7), 384–400. <https://doi.org/10.1108/13527591211281129>
- Aliano, M., Alnabulsi, K., Cestari, G., & Ragni, S. (2023). The role of gender and education in peer-to-peer lending activities: Evidence from a european cross-country study. *European Scientific Journal ESJ*, 2. <https://doi.org/10.19044/esipreprint.2.2023.p95>
- Alistair Milne, & Paul Parboteeah. (2017). *The business models and economics of peer-to-peer lending*. Retrieved from [https://repository.lboro.ac.uk/articles/online\\_resource/The\\_business\\_models\\_and\\_economics\\_of\\_peer-to-peer\\_lending/9494891](https://repository.lboro.ac.uk/articles/online_resource/The_business_models_and_economics_of_peer-to-peer_lending/9494891)
- Amar, M., Ariely, D., Ayal, S., Cryder, C. E., & Rick, S. I. (2011). Winning the battle but losing the war: The psychology of debt management. *Journal of Marketing Research*, 48, S38–S50.
- Ayal, S., Bar-Haim, D., & Ofir, M. (2018). Behavioral biases in peer-to-peer (P2P) lending. *Forthcoming in Behavioral Finance: The Coming of Age (Venezia I. Ed., World Scientific Publishers)*.
- Ayal, S., Hochman, G., & Zakay, D. (2011). Two sides of the same coin: Information processing style and reverse biases. *Judgment and Decision Making*, 6(4), 295–305. <https://doi.org/10.1017/S193029750000190X>
- Ayal, S., & Zakay, D. (2009). The perceived diversity heuristic: The case of pseudodiversity. *Journal of Personality and Social Psychology*, 96(3), 559–573. <https://doi.org/10.1037/a0013906>
- Bondora. (2024, August 28). How are bondora risk ratings calculated? Bondora help center. Retrieved November 13, 2025, from <https://help.bondora.com/hc/en-us/articles/14814705732881-How-are-Bondora-risk-ratings-calculated>
- Cooper, D., Gorbachev, O., & LuengoPrado, M. J. (2023). Consumption, credit, and the missing young. *Journal of Money, Credit and Banking*, 55(2), 379–405.
- Croson, R., & Gneezy, U. (2009). Gender differences in preferences. *Journal of Economic Literature*, 47(2), 448–474.
- Galak, J., Small, D. A., & Stephen, A. T. (2010). Micro-finance decision making: A field study of prosocial lending. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.1634949>
- Herzenstein, M., Dholakia, U. M., & Andrews, R. L. (2011). Strategic herding behavior in peer-to-peer loan auctions. *Journal of Interactive Marketing*, 25(1), 27–36. <https://doi.org/10.1016/j.intmar.2010.07.001>
- Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2), 263–291. <https://doi.org/10.2307/1914185>
- Lee, E., & Lee, B. (2012). Herding behavior in online P2P lending: An empirical investigation. *Electronic Commerce Research and Applications*, 11(5), 495–503. <https://doi.org/10.1016/j.elerap.2012.02.001>
- Liu, Y., Baals, L. J., Osterrieder, J., & Hadji-Misheva, B. (2024). Network centrality and credit risk: A comprehensive analysis of peer-to-peer lending dynamics. *Finance Research Letters*, 63, 105308. <https://doi.org/10.1016/j.frl.2024.105308>
- Manu Siddhartha. (2021). Bondora peer to peer lending loan data. Kaggle. Retrieved November 12, 2025, from <https://www.kaggle.com/datasets/sid321axn/bondora-peer-to-peer-lending-loan-data>
- Ravina, E. (2019). Love & loans: The effect of beauty and personal characteristics in credit markets. *Available at SSRN 1107307*.
- Serrano-Cinca, C., Gutiérrez-Nieto, B., & López-Palacios, L. (2015). Determinants of default in

- P2P lending. *PLOS ONE*, 10(10), e0139427. <https://doi.org/10.1371/journal.pone.0139427>
- Simpson, W. (2001). QAP: The quadratic assignment procedure. *North American STATA Users' Group Meeting*, 1, 1–17.
- Stivala, A., Wang, P., & Lomi, A. (2025). *Improving exponential-family random graph models for bipartite networks*. arXiv. <https://doi.org/10.48550/ARXIV.2502.01892>
- Yao, J., Chen, J., Wei, J., Chen, Y., & Yang, S. (2019). The relationship between soft information in loan titles and online peer-to-peer lending: Evidence from RenRenDai platform. *Electronic Commerce Research*, 19(1), 111–129. <https://doi.org/10.1007/s10660-018-9293-z>
- Zuo, Y. (n.d.). *CLUSTERING ANALYSIS TO SUPPORT LENDER'S DECISION-MAKING IN P2P LENDING*.

## B Technology Statement

During the preparation of this work, we used ChatGPT in order to generate select parts of the R script utilised to process the dataset. Specifically, the tool was used to transform the processed dataset into a format that igraph would accept as a network object. Additionally, some R functions written by the group were improved with ChatGPT's suggestions. No AI tool was utilised to independently write parts of the report. The following parts of the assignment were affected/generated by AI tool usage: **INTRODUCTION**; The tool was utilised to evaluate the validity of certain theoretical concepts and refine them. **DATASET**; the data described within this section was processed partly by some code drafted by ChatGPT and edited by the group. After using this tool/service, **Group 7** evaluated the validity of the tool's outputs, including the sources that generative AI tools have used, and edited the content as needed. As a consequence, **Group 7** takes full responsibility for the content of their work.