# Recommendation Engine for Lower Interest Borrowing on Peer to Peer Lending (P2PL) Platform

Ke Ren
Department of Electrical and Computer Engineering, The
University of Auckland, NZ
keith_kt@foxmail.com

Avinash Malik
Department of Electrical and Computer Engineering, The
University of Auckland, NZ
avinash.malik@auckland.ac.nz

## ABSTRACT

Online Peer to Peer Lending (P2PL) systems connect lenders and borrowers directly, thereby making it convenient to borrow and lend money without intermediaries such as banks. Many recommendation systems have been developed for lenders to achieve higher interest rates and avoid defaulting loans. However, there has not been much research in developing recommendation systems to help borrowers make wise decisions. On P2PL platforms, borrowers can either apply for bidding loans, where the interest rate is determined by lenders bidding on a loan or traditional loans where the P2PL platform determines the interest rate. Different borrower *grades* — determining the credit worthiness of borrowers get different interest rates via these two mechanisms. Hence, it is essential to determine which type of loans borrowers should apply for. In this paper, we build a recommendation system that recommends to any new borrower the type of loan they should apply for. Using our recommendation system, any borrower can achieve lowered interest rates with a higher likelihood of getting funded.

## KEYWORDS

peer-to-peer lending, data mining, recommendation, sentiment analysis

## 1 INTRODUCTION

The development of electronic commerce has lead to a burgeoning growth in online Peer to Peer Lending (P2PL) system. P2PL system is a micro financing platform, which is rising as an alternative to traditional financial lenders such as banks. There are two main participants in P2PL systems: borrowers and lenders. On the one side, borrowers apply for loans. On the other side, lenders can view the characteristics of the borrowers/loans and decide, which loans they should invest in. In recent years, a great deal of research has gone into developing recommendation systems to help lenders [7, 15] achieve high returns with low risk of defaults. However, there has

**Table 1: Average interest rates of traditional loans and bidding loans for borrowers with the same characteristics.**

| Grade | AA | A | B | C | D | E | HR |
|---|---|---|---|---|---|---|---|
| Average traditional interest | 0.112 | 0.082 | 0.158 | 0.197 | 0.247 | 0.295 | 0.318 |
| Average bidding interest | 0.113 | 0.102 | 0.151 | 0.182 | 0.208 | 0.247 | 0.235 |
| Traditional − bidding | -0.001 | -0.02 | 0.007 | 0.015 | 0.039 | 0.048 | 0.083 |

**Table 2: T-test between the interest rates of traditional loans and bidding loans for each grade, with null hypothesis that they have the same mean value.**

| Grade | AA | A | B | C | D | E | HR |
|---|---|---|---|---|---|---|---|
| P-value | 0.74 | 3.6e-09 | 0.061 | 1.41e-05 | 1.15e-22 | 1.40e-27 | 1.77e-29 |
| Decision | Not reject | Reject | Not reject | Reject | Reject | Reject | Reject |

not been much research into developing recommendation systems to advice borrowers. In particular, the main objective from borrower's perspective is getting funded with the lowest interest rate payable. We build a recommendation framework for borrowers to help them borrow with lower interest rates and increased likelihood of getting funded on P2PL platforms in this paper.

From the borrower's perspective, there are two essential questions that need to be considered when applying for loans: ① will the loan be funded successfully? ② What is the lowest obtainable interest rate? Online P2PL platforms do not help borrowers with these two questions, but rather give the borrowers a choice to select from different types of loans that they can apply for. On online P2PL platforms[1], the two main types of loans are:

- **Traditional loan:** based on the borrower's personal information, P2PL platforms decide the interest rate for each borrower's loan. Next, the P2PL platforms put the loan online for a certain period for lenders to fund the loan.
- **Bidding loan:** first and foremost, borrowers themselves decide the maximum interest rate they are willing to pay. Then P2PL platforms put the loan online and wait for lenders to bid on the loan, with the interest rate that they want. At the end of the bidding period, if the loan receives sufficient funding, P2PL platforms will select lenders with the lowest interest rate. However, if the final total interest rate is higher than the borrower's maximum selected interest rate, then this loan is *not* funded.

To show that it is necessary to select the right type of loan when applying on P2PL platforms, consider the average historical interest rate for traditional and bidding loans from Prosper (one of the largest P2PL platform in the world) in Table 1. Table 1 shows the interest rates of traditional and bidding loans for each borrower grade along with their differences. A higher grade (e.g., AA) indicates lower likelihood of the borrower defaulting and a lower grade

---

[1]prosper.com

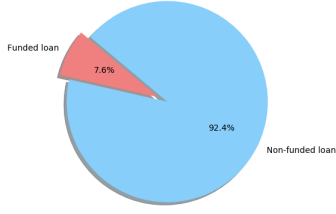Figure 1: Pie chart of distribution on 12006 bidding loans

| AA | A | B | C | D | E | HR |
|-------|-------|-------|-------|-------|------|------|
| 34.2% | 33.1% | 27.3% | 16.1% | 10.4% | 4.7% | 1.6% |

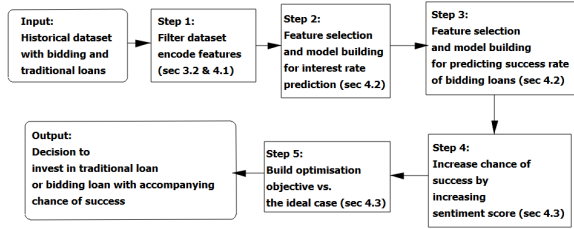Table 3: Average success rate of funding bidding loans of different grades.



Figure 2: Overall proposed methodology

(e.g., HR) indicates higher likelihood of the borrower defaulting on their loan obligations. The T-test with the null-hypothesis that the traditional and bidding loans have the same mean value is shown in Table 2. We can observe from Tables 1 and 2 that borrowers with credit grade A should apply for a traditional loan, while borrowers with lower credit grades C, D, E, and HR would achieve a lower interest rate payable when applying for bidding loans. Finally, borrowers with credit grade AA and B can either apply for bidding or traditional loans, since there is no significant statistical difference between the interest rates of bidding and traditional loans, for these grades. Especially for borrowers with HR grade, the interest rate payable, when applying for a bidding loan, is decreased by 8.3%. Hence, it is necessary for borrowers to decide, which types of loan should they apply for.

Getting a lower interest rate is one borrower objective, the other objective is to actually get funded. Figure 1 illustrates the distribution of funded and non-funded bidding loans from a total of 12006 loans from the Prosper historical dataset. It can be seen that on *average* only 7.6% of all bidding loans get funded. The success of getting funded for different grades of bidding loans is shown in Table 3. Thereby making it important for borrowers to make a wise choice when applying for a loan.

Our **major contribution** in this work is to build a recommendation system for borrowers on P2PL platform, which takes as input the historical loan data with the borrower's characteristic and outputs the decision on the types of loans they should apply for. Using our recommendation system, borrowers can achieve a reduced interest rate payable, with a higher chance of successfully getting the loan request funded. The overview of our proposed technique and key technical contributions are shown in Figure 2.

1. We start with the historical loan dataset. We first filter the dataset to remove unusable rows. Next, we encode the categorical features to numerical ones.
2. We build machine learning models to predict the interest rate payable for bidding and traditional loans along with selecting the most important borrower features that influence the models.
3. We build machine learning models to classify if a given borrower will succeed on the bidding loan platform along with selecting the features that positively influence the model. In steps ② and ③, we compare different machine learning algorithms, including linear and logistic regression (LOGIT) [13], Random Forest (RF) [6], Support Vector Machine (SVM) [5], and k-nearest neighbour (k-NN) [4].
4. We improve the positive sentiments in the textual description for borrowing, which in turn increases chances of bidding loan being successfully funded.
5. Using the results from the previous steps, we compare the interest rate and success rate of traditional and bidding loans with the ideal case: 0% interest rate and 100% success rate. The one closest to the ideal case would be recommended as the loan type that the borrower should apply for.

The rest of the paper is organised as follows. Section 2 reviews and discusses the current state-of-the-art. Section 3 describes the problem statement and the P2PL dataset. Section 4 describes the details of the workflow of our proposed technique. The experimental results and quantitative comparison with the current state-of-the-art technique is presented in Section 5. Finally, we conclude the paper and discuss the advantages and limitations of the proposed model in Section 6.

## 2 RELATED WORK

With the burgeoning growth of online P2PL marketplaces, a great deal of research has been proposed to guide lenders and borrowers to benefit from the P2PL system. From the lender's perspective, recent work in [12] compares different machine learning algorithms and finds that the best algorithm to predict the possibility of a loan/borrower defaulting is random forests. Another work in [9] studies the strategic herding behaviour in P2PL loan auctions and points out that the strategic herding behavior benefits bidders individually and collectively. Other works in [7] and [15] proposed recommendation systems for lenders that yield an investment portfolio with minimal risk of default along with maximum returns.

From a borrower's perspective, recent works in [10] and [1] study the role of identity claims constructed in narratives by borrowers, and reveals that as the number of identity claims in narratives increases, the likelihood of successful funding also increases. Another work in [8] studies the determinants of funding success in online P2PL communities and finds out that the most predominant predictors of loan's likelihood of being funded successfully are the extent of personal characteristics[2] provided by borrowers, and their credit grades. The work in [16] weighs the financial and social features of borrowers to determine their influence in the success of loan being funded. For the purpose of predicting the likelihood of the loan being funded successfully, the most recent work in [3] explores

---

[2] The text describing the reason for borrowing money.

temporal dynamics of loan listings and builds a regression model to predict likelihood of successful funding. However, this model can only yield high accuracy under the assumption that the bidding process on loans is already finished. Specifically, the work in [3] uses features called 'number of bids', that are recorded in historical dataset *only* after the bidding process is completed. We aim to help **new** borrowers to make good decisions on the types of loans to apply for. In turn, this means that features recorded after completion of bidding cannot be used in our problem setup. In this paper, we compare different machine learning algorithms with three feature selection techniques to improve the accuracy of prediction. Moreover, we also quantitatively compare our technique with the one proposed in [3].

## 3 PROBLEM STATEMENT, OBJECTIVES, AND DATA ANALYSIS

In the next two sections we give our problem statement, objectives, and describe the historical dataset we have used for analysis and prediction.

### 3.1 The problem statement and objectives

In a P2PL system, let: ① $U = \{u_1, u_2, \ldots, u_n\}$ be the set of borrowers, ② $C = \{c_1, c_2, \ldots, c_m\}$ be the set of features (characteristic of the borrowers) and ③ $V = \{v_{ij} : i \in U, j \in C\}$ be the set of values for each feature of each borrower, where the entries $v_{ij}$ denote the value of feature $j$ given by borrower $i$. Furthermore, we use $P_{bid} \subseteq C$ and $P_{trad} \subseteq C$ as the set of predictors for bidding loans and traditional loans needed to predict the interest rates $I_{bid}$ and $I_{trad}$, respectively. In practice, the traditional loans are funded quickly. The work in [2] finds the success rate of getting funded for traditional loans is 81%. Thus, we set the success rate of the traditional loan $S_{trad} = 0.81$. Finally, we use $P_{suc} \subseteq C$ as the set of predictors for success rate of bidding loans $S_{bid}$. Overall, we specially focus on the following research problems:

- How to accurately predict the interest rate of bidding and traditional loans, $I_{bid}$ and $I_{trad}$, respectively.
- How to accurately predict and increase the success rate of getting funded for bidding loans, $S_{bid}$.
- Given the interest rates and success rates of both bidding and traditional loans. In order to obtain the lowest interest rate payable and increase his/her chances of successfully getting funded, which type of loans should the borrower apply for? This can be formulated as below:

$$\max \quad (-I, S)$$
$$\text{s.t.} \quad (I, S) \in \{(I_{trad}, S_{trad}), (I_{bid}, S_{bid})\},$$
$$\text{and} \quad I, S \in [0, 1],$$

where $I, S \in \mathbb{R}^{\geq 0}$.

### 3.2 Data analysis

In this paper, we train and test based on dataset from a well established online P2PL platform, Prosper. We choose Prosper, because Prosper is America's first online P2PL platform with over $14 billion in funded loans since 2006 [14]. In addition, Prosper has been in operation for more than 10 years, and hence, can offer a plethora of historical data that is necessary for training and testing. Other

researchers who study P2PL systems also use the same dataset from Propser [17], which makes it possible for us to compare our work with the current state-of-the-art techniques.

In this paper, we use *two* Prosper datasets. ① The *traditional loan* dataset, which contains 113,938 funded loans with 81 features in total, dating from 2005 to 2014. ② The *bidding loan* dataset, which contains 12,774 bidding loans with 30 features in total, dating from 2007-5-27 to 2007-6-30. Not all features, from both dataset, are applicable to our study, hence we filter the dataset using the following rules:

- Remove blank, zero, and missing features, because there is no information in these features.
- Remove features that are *not* applicable to new borrowers, for example 'number of bids' and 'loan current days delinquent'. We are aiming to predict the interest rate payable and success rate of loan getting funded for *new* borrowers, any feature that are recorded after the loan has started cannot be considered.
- Remove loans with missing values.

After filtering both datasets by applying the above rules, we end up with two datasets that are described below:

- *Traditional dataset*: contains 70,849 funded loans with 31 features and 1 response variable — the borrower's interest rate. Among these 31 features, 5 of them are categorical and the rest are numerical (see Table 4).
- *Bidding dataset*: contains 12,006 loans (both funded and non-funded loans) with 12 features and 2 response variables — the borrower's interest rate and the status of the bidding loan — funded or not funded. Among these 12 features, 6 of them are categorical and the rest are numerical (see Table 5).

**Table 4: Features and response variable description of traditional dataset**

| Feature | Explanation | Type |
| --- | --- | --- |
| BorrowerRate | The Borrower's interest rate for this loan. | Numerical |
| OpenCreditLines | Number of open credit lines. | Numerical |
| ProsperGrade | A custom rating score built by Prosper. | Categorical |
| ProsperScore | A custom risk score built by Prosper. | Numerical |
| ListingCategory | The category of the listing that the borrower. | Numerical |
| CurrentCreditLines | Number of current credit lines. | Numerical |
| TotalCreditLinespast7years | Number of credit lines in the past seven years. | Numerical |
| OpenRevolvingAccounts | Number of open revolving accounts. | Numerical |
| OpenRevolvingMonthlyPayment | Monthly payment on revolving accounts. | Numerical |
| TotalInquiries | Total number of inquiries. | Numerical |
| CurrentDelinquencies | Number of accounts delinquent. | Numerical |
| AmountDelinquent | Dollars delinquent. | Numerical |
| Occupation | The Occupation selected by the Borrower. | Categorical |
| PublicRecordsLast10Years | Number of public records in the past 10 years. | Numerical |
| RevolvingCreditBalance | Dollars of revolving credit. | Numerical |
| TradesNeverDelinquent | Trades that have never been delinquent. | Numerical |
| TotalTrades | Number of trade lines ever opened. | Numerical |
| StatedMonthlyIncome | The monthly income the borrower stated. | Numerical |
| AvailableBankcardCredit | The total available credit via bank card. | Numerical |
| TradesOpenedLast6Months | Number of trades opened in the last 6 months. | Numerical |
| BankcardUtilization | The percentage of available credit that is utilized. | Numerical |
| Homeownership | Specifies if the borrower is a homeowner or not. | Categorical |
| DebtToIncomeRatio | The debt to income ratio of the borrower. | Numerical |
| InquiriesLast6Months | Number of inquiries in the past six months. | Numerical |
| LoanAmount | The origination amount of the loan. | Numerical |
| CreditScoreRangeLower | The lower range of the borrower's credit score. | Numerical |
| EmploymentStatusDuration | The length in months of the employment status. | Numerical |
| DelinquenciesLast7Years | Number of delinquencies in the past 7 years. | Numerical |
| Term | The length of the loan expressed in months. | Numerical |
| BorrowerState | The state of the address of the borrower. | Categorical |
| EmploymentStatus | The employment status of the borrower. | Numerical |
| Description | The description of the lowan written by the borrower. | Categorical |

**Table 5: Features and response variable description of bidding dataset**

| Feature | Explanation | Type |
|---------|-------------|------|
| BorrowerRate | The Borrower's interest rate for this loan. | Numerical |
| BorrowerMaximumRate | The maximum interest rate the borrower will accept. | Numerical |
| ProsperGrade | A custom rating score built by Prosper. | Categorical |
| Homeownership | Specifies if the borrower is a homeowner or not. | Categorical |
| DebtToIncomeRatio | The debt to income ratio of the borrower. | Numerical |
| LoanAmount | The origination amount of the loan. | Numerical |
| FundingOption | The options of funding. | Categorical |
| Images | Number of images that are uploaded by borrowers. | Numerical |
| Duration | The length of funding duration. | Numerical |
| BorrowerState | The state of the address of the borrower. | Categorical |
| EmploymentStatus | The employment status of the borrower. | Numerical |
| HasVerifiedBankAccount | Specifies if or not the bank account is verified. | Categorical |
| Description | The description of the lowan written by the borrower. | Categorical |
| LoanStatus | The current status of the loan. | Categorical |

**Table 6: A simple example: characteristic of two borrowers**

| Feature | Borrower 1 | Borrower 2 |
|---------|-----------|-----------|
| Borrower maximum rate | 0.16 | 0.105 |
| Prosper grade | 7 (HR) | 1 (AA) |
| Term | 36 | 36 |
| Credit score | 540 | 760 |
| Delinquencies in last 7 years | 5 | 0 |
| Debt to income ratio | 0.17 | 0.06 |
| Loan amount | 2,300 | 10,000 |
| Homeownership | 0 (Not own) | 1 (Own) |
| Duration | 3 | 10 |
| Funding option | 0 (Close when funded) | 1 (Open for duration) |
| Has verified bank account | 1 (True) | 1 (True) |
| Images | 0 | 0 |
| Description | 0.3818 (Payoff Credit Cards) | 0 (Lender seeing Prosper from borrower's point-of-view) |

## 4 METHODOLOGY

In this section, we first introduce the method we use to encode categorical features followed by sentiment analysis for textual descriptions of reasons for borrowing, as input by the borrowers. Next, we introduce the feature selection techniques and machine learning algorithms we use to predict interest rate payable for bidding and traditional loans, respectively, along with computing the chances of success of any given bidding loan (recall that the chance of success of traditional loan is fixed at 81%). Finally, we propose the method to advice borrowers with the type of loans they should apply for.

### 4.1 Feature encoding and sentiment analysis

Both bidding and traditional dataset have several categorical features. These features need to be transferred to a numerical value, so that they can be used in machine learning algorithms like linear regression, Logistic Regression (LOGIT), etc. In the biding and traditional datasets, we split the categorical features into three types:

1. Features that only have two classes such as "Homeownership" and "Funding option".
2. Features that have more than two classes like "Prosper grade" and "Borrower's state".
3. Features that contain random textual (English) words like Borrower's "Description" of the loan.

Table 6 details the features of interest for two borrowers. In order to encode the categorical features, we use two most popular encoding techniques: *binary encoding* and *ordinal encoding*. Binary encoding technique transfers each categorical feature into *new* numerical features containing only zeros and ones. For instance, the categorical feature "Homeownership" has two classes "Own" or "Not own". Each of these classes is encoded as an individual numerical feature, with a value of 0 or 1. In case of the first borrower in Table 6, the "Homeownership" feature will be translated into

two new features "Not own:1" and "Own:0". However, either of these two new features is enough to show the home ownership of the borrower. To avoid the increase of the number of features, we only select one of them as the encoded feature. In other words, the feature "Homeownership" is encoded with a value of 0 (Not own) or 1 (Own). Same for borrower two.

Binary encoding is only feasible for categorical features with few classes. In case of a categorical feature with a plethora of classes, this method will increase the total number of features in the dataset significantly. Ordinal encoding is the preferred method in such cases. It converts string labels to integer values 1 through $k$, where $k$ is the number of classes in a given categorical feature. For example, consider the feature "Proser grade" shown in Table 6. There are 7 classes in this feature: AA, A, B, C, D, E, HR. After applying ordinal encoding, the 7 classes are transferred to 1, 2, 3, 4, 5, 6 and 7, respectively. Finally, neither binary encoding nor ordinal encoding is applicable when encoding the third type of categorical feature: textual data, because textual description has meaning, which should be captured by the encoding technique.

In order to encode the third type of feature, we do sentiment analysis. Sentiment analysis can extract and evaluate the emotions in text. There are two popular types of sentiment analysis: ① classify the polarity of given text as positive, negative or neutral. ② Evaluate a given piece of text to a certain score, which shows the levels of emotion. In this paper, we apply the second type of sentiment analysis to encode the text as numerical value. Specifically, we apply the sentiment analysis technique from VADER (Valence Aware Dictionary and sEntiment Reasoner) [11], thereby encoding the text into numerical scores, which we call the *sentiment score*, ranging from -1 to 1. Here, 1 represents the most positive emotion and -1 repents the most negative emotion. For example, in Table 6, the description "Payoff Credit Cards" (as the reason for borrowing money) is evaluated to 0.3818. An optimal sentiment score can help with getting the loan funded. The results comparing the likelihood of getting funded with varying sentiment scores are described in Section 5.4.

### 4.2 Machine learning models for interest rate prediction, likelihood of getting funded, and feature selection

Recall that there are 31 features in the traditional dataset and 12 features in the bidding dataset (Section 3.2). However, not all of these features are useful when predicting the interest rate and/or the success rate of getting funded. Moreover, the machine learning algorithms we use in this paper such as SVM and k-NN are sensitive to irrelevant features. Hence, it is necessary to do feature selection along with predicting the interest rates payable and the likelihood of getting funded.

In this paper, we compare three popular feature selection algorithms: ① forward selection, ② backward selection and ③ recursive selection. Since all the three techniques are well known, here we only give an overview of the techniques and the results are shown in Sections 5.2 and 5.3. In forward selection, we start with an empty feature set ($P_{trad}$, $P_{bid}$, and $P_{suc}$) and keep on adding new features one after another until the coefficient of determination/recall rate stops increasing. In backward selection, we start

all the features in the feature sets, and keep on dropping features one after another until the coefficient of determination/recall rate stops increasing. In recursive selection, we start with a set of all the features in the feature sets and keep on dropping the least important features[3] until the coefficient of determination/accuracy stops increasing.

In order to predict the interest rates of bidding and traditional loans, we choose four regression models to compare and select the best fitted model. These regression models include linear regression, Random Forest (RF), Support Vector Machine (SVM) and $k$-Nearest Neighbors (k-NN). All four models have both advantages and disadvantages, and we apply feature selection on each of the model to find the best fit. Different from predicting the interest rates, we select four machine learning classifiers to predict the likelihood of a bidding loan getting funded. RF, SVM, and k-NN are all applicable for classification and regression problem, we only replace linear regression with LOGIT to be the fourth classifier. Again, we apply feature selection on each of the classifier and find the classifier that results in the highest accuracy (recall rate using cross validation techniques). The results are described in Section 5.

### 4.3 The decision process to recommend the type of loan application

Our final goal is to help new borrowers decide, which type of loan they should apply for. The goal is to achieve the highest likelihood of successfully getting funded at the lowest interest rate payable. To reach this goal, we first compute the interest rate payable and the likelihood of success using the models described in Section 4.2. The machine learning models output two tuples: $(I_{trad}, 0.81)$ and $(I_{bid}, S_{bid})$. We next compare these two tuples with the ideal case: $(0, 1)$, where the first element of the tuple indicates 0% interest rate payable and the second indicates 100% likelihood of successfully getting funded. The final decision is made by comparing the Euclidean distance between each of the tuples obtained from the machine learning algorithms and the ideal case. We can formalise the approach as follows:

$$
\begin{aligned}
\min \quad & |(I, S) - (I_{ideal}, S_{ideal})| \\
\text{s.t.} \quad & (I_{ideal}, S_{ideal}) = (0, 1), \\
& (I, S) \in \{(I_{trad}, S_{trad}), (I_{bid}, S_{bid})\}, \\
\text{and} \quad & I, S \in [0, 1],
\end{aligned}
\tag{1}
$$

where $|\cdot|$ represents the Euclidean distance, and $I, S \in \mathbb{R}^{\geq 0}$.

## 5 EXPERIMENTAL RESULTS

In this section, a thorough comparison of the various techniques described in Section 4 is presented.

### 5.1 Experimental setup

All our experiments are performed on the traditional and bidding Prosper datasets obtained from [17] and cleaned/analysed as described in Section 3.2. The detail of datasets and experiments performed is presented below:

---

[3]Least important features are obtained after fitting the model using the scikit learn Python library.



**Figure 3: Prediction of interest rates for traditional loans.**

- To predict the interest rate of traditional loans we randomly sample 10,000 loans from the 70,849 loans available to us. For training and testing, we perform 5-fold Montecarlo cross validation by splitting the 10,000 loans into a ratio of 80:20. The results of the 5 runs and the average are shown in Section 5.2.
- To predict the interest rate of bidding loans we sample 908 funded bidding loans from the 12,006 available in the bidding dataset. For training and testing, we perform 5-fold Montecarlo cross validation by splitting the 908 loans into a ratio of 80:20 for training and testing, respectively. The results of the 5 separate runs and the average results are shown in Section 5.2.
- To predict the success rate of bidding loans we sample 908 funded and 908 non-funded loans from the bidding dataset to get in total 1816 loans. Again, for training and testing, we perform 5 fold Montecarlo cross validation, with a 80:20 split for each run, the results are presented in Section 5.3.
- The results of sentiment scores impacting success rate of getting funded and the overall efficacy of the recommendation system are presented in Sections 5.4 and 5.5.

### 5.2 Interest rate payable prediction for traditional and bidding loans

In this section, we compare four machine learning regression models with feature selection techniques to predict the interest rates payable for both traditional and bidding loans. Since the response variable, interest rate, is continuous, the coefficient of determination $R^2$ is used to evaluate the accuracy of models for comparison purposes. The coefficient of determination can be formulated as:

$$
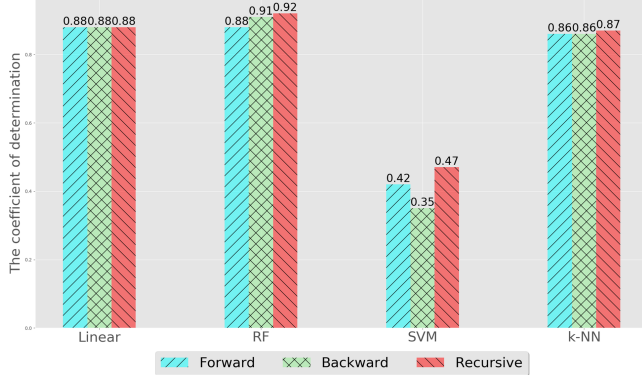R^2 = \frac{\text{Residual sum of squares}}{\text{Total sum of squares}}
$$

where $R^2$ ranges from 0 to 1. The larger the value of $R^2$, the better the goodness of fit of a model. Specifically, an $R^2$ value equal to 1 means the regression model perfectly fits the data.

Figure 3 illustrates the average performance of predicting the interest rates payable, for 5 cross validation runs, for traditional loans by applying linear regression, RF, SVM and k-NN with forward, backward and recursive feature selection. From Figure 3, we can observe that under the same feature selection technique, RF always performs the best. This result also can be seen from Table 7, which shows the 5 Monte-Carlo cross validation runs for each model with

**Table 7: The 5 Monte-Carlo CV of the four regression models with recursive selection for predicting the interest rate of traditional loans**

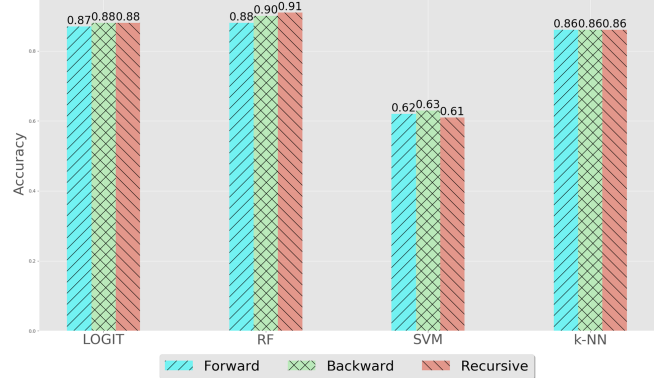| CV test | 1 | 2 | 3 | 4 | 5 | Average |
|---|---|---|---|---|---|---|
| Linear | 0.91 | 0.94 | 0.94 | 0.93 | 0.92 | 0.93 |
| RF | 0.93 | 0.97 | 0.96 | 0.97 | 0.96 | 0.96 |
| SVM | 0.57 | 0.62 | 0.60 | 0.59 | 0.61 | 0.60 |
| k-NN | 0.90 | 0.94 | 0.94 | 0.93 | 0.95 | 0.93 |



**Figure 4: Prediction of interest rates for bidding loans.**

**Table 8: The 5 Monte-Carlo CV of the four regression models with recursive selection for predicting the interest rate of bidding loans**

| CV test | 1 | 2 | 3 | 4 | 5 | Average |
|---|---|---|---|---|---|---|
| Linear | 0.89 | 0.88 | 0.91 | 0.87 | 0.86 | 0.88 |
| RF | 0.92 | 0.93 | 0.92 | 0.93 | 0.90 | 0.92 |
| SVM | 0.47 | 0.45 | 0.46 | 0.49 | 0.48 | 0.47 |
| k-NN | 0.89 | 0.87 | 0.89 | 0.85 | 0.84 | 0.87 |

recursive feature selection. Recursive feature selection outperforms forward and backward feature selection under the same regression model. In addition, RF with recursive selection achieves the highest coefficient of determination with value 0.96. Linear regression and k-NN also fit the data well with recursive selection with $R^2 = 0.93$. However, SVM does not perform well with the value of $R^2$ of only 0.6. We choose RF with recursive selection to be the best method to predict the interest rate of traditional loans. The features selected by recursive feature selection technique with RF are: Prosper grade, term, credit score and delinquencies in last 7 years.

Next, we predict the interest rates of bidding loans. Figure 4 illustrates the average performance of predicting the interest rates payable, for 5 cross validation runs, for bidding loans by applying linear regression, RF, SVM and k-NN with forward, backward and recursive feature selection. It can be observed from Figure 4 that under the same feature selection technique, RF performs best. The highest $R^2$ of 0.92 is achieved by applying RF with recursive feature selection. Linear regression and k-NN also perform well with the values of $R^2$ of 0.88 and 0.87, respectively. In addition, with the same regression model, recursive feature selection gives the best subset of features that results in the highest value of $R^2$. This result also can be observed from Table 8, which describes the 5 Montecarlo



**Figure 5: Prediction of the likelihood of successfully getting funded for bidding loans.**

**Table 9: The 5 Monte-Carlo CV of the four regression models with recursive selection for predicting the success of getting funded of bidding loans**

| CV test | 1 | 2 | 3 | 4 | 5 | Average |
|---|---|---|---|---|---|---|
| LOGIT | 0.86 | 0.87 | 0.89 | 0.88 | 0.90 | 0.88 |
| RF | 0.89 | 0.93 | 0.90 | 0.92 | 0.92 | 0.91 |
| SVM | 0.59 | 0.61 | 0.63 | 0.60 | 0.61 | 0.61 |
| k-NN | 0.87 | 0.85 | 0.89 | 0.87 | 0.88 | 0.87 |

cross validation runs with recursive feature selection. SVM again does not perform well when predicting the interest rate payable for bidding loans with the value of $R^2$ of only 0.47. Therefore, for the purpose of accurately predicting the interest rate of bidding loans, we select RF with recursive feature selection as the preferred prediction model. The selected features are: borrower maximum rate, Prosper grade, debt to income ratio, loan amount, homeownership, duration, funding option and has verified bank account.

## 5.3 Predicting the success rate of funding bidding loans

In this section, we compare LOGIT, RF, SVM, and k-NN with forward, backward and recursive feature selection techniques to find the best classification model for predicting the success rates of getting funded for bidding loans. Since the response variable here is either funded or non-funded, we select the accuracy (recall rate) as the criterion to evaluate the goodness of fit of the model. The accuracy measure used for comparison purpose can be formulated as follows:

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of predictions}}$$

where the accuracy $\in [0, 1]$. The higher the accuracy, the better the model fits the dataset.

Figure 5 shows the average accuracy, from amongst the 5 cross validation runs, of predicting the success of getting funded for bidding loans by applying LOGIT, RF, SVM and k-NN with forward, backward and recursive feature selection. The results of the 5 Monte-Carlo cross validation runs with recursive feature selection are shown in Table 9. From Figure 5, it can be seen that with the same feature selection technique, RF has the best average accuracy. In

**Table 10: The confusion matrix of RF with recursive selection on test bidding dataset**

|  | Actual funded | Actual non-funded |
|---|---|---|
| Predicted funded | 161 true positives | 20 false positives |
| Predicted non-funded | 12 false negatives | 171 true negatives |

**Table 11: The accuracy of predicting the success of getting funded of bidding loans**

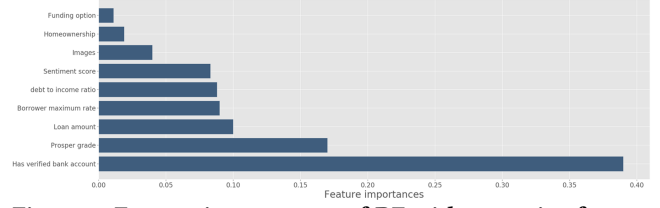| Algorithm | Accuracy | Selected features |
|---|---|---|
| RF with recursive selection | **0.91** | Borrower maximum rate, Prosper grade, debt to income ratio, loan amount, homeownership, funding option, has verified bank account, images, sentiment score |
| Current state-of-the-art | 0.67 | Borrower maximum rate, debt to income ratio, loan amount, homeownership, listing description length |

addition, the performance of SVM is bad around 62% accuracy. Therefore, to predict the success of getting funded of bidding loans as accurately as possible, we select the RF together with recursive feature selection as the preferred prediction model.

Table 10 shows the confusion matrix of RF with recursive feature selection for a single Montecarlo run . We can observe that the false positives and false negatives are small compared to the true positives and the true negatives. Specifically, the true positive rate is 93% and the true negative rate is 90%. These results validate that our proposed model is reliable when predicting the success of getting funded for bidding loans.
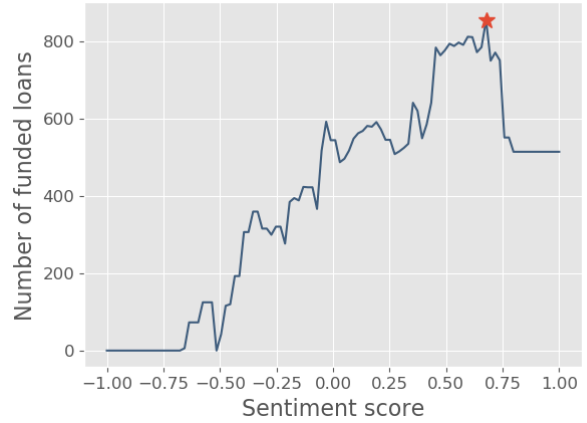
In order to verify that the proposed technique is better than the current state-of-art model [3], we randomly split the bidding dataset 1816 loans (908 funded and non-funded loans) to training and testing dataset with ratio of 80:20. Then we train the RF model with recursive feature selection and the current state-of-the-art model on the training dataset (1452 loans) and test on the testing dataset (364 loans). The results are shown in Table 11. We can observe from Table 11 that our proposed model has 0.91 accuracy, which is 0.24 (24%) higher than the current state-of-the-art technique. Hence, we can state that the RF technique with the recursive feature selection algorithm outperforms the current state-of-the-art technique (as proposed in [3]).

### 5.4 Impact of sentiment score

In this section we apply the sentiment analysis technique from VADER [11] on the bidding dataset and study the impact of raising the sentiment score on the likelihood of getting funded for bidding loans. Recall that RF with recursive feature selection is the best model to predict the success rate of getting funded for bidding loans. The feature "sentiment score" is selected in the predictors (see Table 11). Recall that the feature sentiment score is encoded by applying sentiment analysis from VADER on the textual "description" for borrowing. The selection of sentiment score, as a predictor, indicates that the emotion of texts written by borrowers does impact the success rate of getting funded when applying for

**Figure 6: Feature importances of RF with recursive feature selection model**

**Figure 7: Impact on success of getting funded by changing the sentiment of texts.**

**Table 12: Results of increasing sentiment score on the 12,006 loans in the bidding dataset.**

| Dataset | Funded loan | Non-funded loan |
|---|---|---|
| No change in sentiment score | 908 | 11098 |
| Increasing the sentiment score to a positive value | 1764 | 10242 |

bidding loans. Figure 6 gives the feature importances of the nine features selected by RF with recursive feature selection technique. From Figure 6, we can observe that sentiment score is the $6^{th}$ most influential feature.

We want to observe if crafting the "description" for borrowing can help improve the overall success rate of getting funded. In order to do so, we control the sentiment score of all 11,098 non-funded loans, from the bidding dataset, from -1 to 1 and then apply the proposed success rate prediction model to see how many of them will become funded. Figure 7 shows the number of funded loans by changing the sentiment score. We can observe that by changing the sentiment score to anywhere between (0.45, 0.70), around 800 non-funded loans are transferred to funded. In addition, when sentiment score equals 0.68, the maximum number of loans (856) get funded. These results indicate that borrower' should write the description with a more positive sentiment. However if the description is *too* positive, the chance of getting funded decreases, which can be seen from Figure 7. This is because too much positive sentiment looks fake.

Table 12 compares the number of originally funded loans in the bidding dataset and the new dataset obtained by changing the

**Table 13: Comparison between the historical data and the results obtained by applying the proposed mehtod**

| | Traditional loans | Bidding loans | Funded | Non-funded | Average interest rate of funded loans |
|---|---|---|---|---|---|
| Historical/original dataset | 500 | 500 | 569 | 431 | 0.23 |
| The proposed method | 892 | 108 | 820 | 180 | 0.20 |

sentiment score to a positive value. It can be observed that the number of funded loans in the resultant dataset, which is obtained by increasing the sentiment score, is about twice that in the original dataset. Performing a t-test with the null hypothesis that these two dataset have the same mean, the p-value is just 1.67e-69. Hence, we can deduce that our proposed method can potentially improve the chance of success of getting funded significantly.

## 5.5 Efficacy of the overall recommendation engine

In this section we present the results of the overall efficacy of the proposed recommendation engine. We first sample 1000 loans from the historical dataset; 500 from the traditional and 500 from the bidding dataset, respectively. Next, we apply RF with recursive feature selection to obtain the interest rates payable ($I_{trad}$, $I_{bid}$) for each loan. Next, we compute the success rate of the loans ($S_{bid}$) after setting the sentiment score to 0.68 (the optimal value), recall that $S_{trad} = 0.81$. Finally, we apply Equation (1) to recommend to the borrower if they should be applying for traditional or bidding loan.

The results are shown in Table 13. From Table 13, we can observe that most of the borrowers who apply for the bidding loans are recommended to apply for traditional loans. This results is expected, because most biding loans, even after increasing the sentiment scores, remain unfunded. However, the number of funded loans increases from 569 to 820, and the average interest rate of funded loans is decreased by 3%. These results show that our proposed recommendation system can help borrowers to get fund successfully with a lower interest rate.

## 6 CONCLUSION AND FUTURE WORK

Online Peer to Peer Lending (P2PL) marketplaces, connect lenders directly to borrowers. The convince of use has led to a burgeoning growth of P2PL marketplaces. Borrowers using P2PL marketplaces are usually looking to get a loan at the lowest interest rate. P2PL marketplaces provide two main pathways for obtaining loans: ① traditional loans, where the platform decides the interest rate for the borrowers and the lenders fund the loan, or ② the lenders themselves decide on the interest they want, by bidding on the loan. Borrowers can get a lowered interest rate via the bidding technique, however with a reduced likelihood of getting funded. Hence, it is essential to make individual recommendation to borrowers depending upon their situation. However, to the best of our knowledge no such recommendation system exists.

We build a recommendation system for borrowers that recommends the best loan option for them, which results in higher likelihood of getting funded, while reducing the interest rate payable. Our methodology consists of three main steps. In step-①, use RF with recursive borrower feature selection model to predict the interest rates of bidding and traditional loans. In step-②, we use RF with recursive borrower feature selection model to predict the success rate of getting funded for bidding loans and even improve

the sentiment of reasons for borrowing. Finally, in step-③, given the interest rates and success rates of both traditional and bidding loans, we compare them with the ideal case and determine the best choice for borrowers.

Experimental results show that our proposed method outperforms the current state-of-the-art technique, in that the accuracy of correctly predicting the success rate of bidding loans increases from 67% to 91%. In addition, the proposed technique can increase the chances of getting funded by 2×. The main drawback of our proposed method is that it currently cannot craft the reason for borrowing in order to increase the sentiment scores, we plan to remedy this in the future.

## REFERENCES

[1] Tim Althoff, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. 2014. How to ask for a favor: A case study on the success of altruistic requests. In *Eighth International AAAI Conference on Weblogs and Social Media*.

[2] Nataliya Barasinska and Dorothea Schäfer. 2014. Is crowdfunding different? Evidence on the relation between gender and funding success from a German peer-to-peer lending platform. *German Economic Review* 15, 4 (2014), 436–452.

[3] Simla Ceyhan, Xiaolin Shi, and Jure Leskovec. 2011. Dynamics of bidding in a P2P lending service: effects of herding and predicting loan success. In *Proceedings of the 20th international conference on World wide web*. ACM, 547–556.

[4] Samprit Chatterjee and Seymour Barcun. 1970. A nonparametric approach to credit screening. *Journal of the American statistical Association* 65, 329 (1970), 150–154.

[5] Harris Drucker, Christopher JC Burges, Linda Kaufman, Alex J Smola, and Vladimir Vapnik. 1997. Support vector regression machines. In *Advances in neural information processing systems*. 155–161.

[6] David Feldman and Shulamith Gross. 2005. Mortgage default: classification trees analysis. *The Journal of Real Estate Finance and Economics* 30, 4 (2005), 369–396.

[7] Yanhong Guo, Wenjun Zhou, Chunyu Luo, Chuanren Liu, and Hui Xiong. 2016. Instance-based credit risk assessment for investment decisions in P2P lending. *European Journal of Operational Research* 249, 2 (2016), 417 – 426. https://doi.org/10.1016/j.ejor.2015.05.050

[8] Michal Herzenstein, Rick L Andrews, Utpal M Dholakia, and Evgeny Lyandres. 2008. The democratization of personal consumer loans? Determinants of success in online peer-to-peer lending communities. *Boston University School of Management Research Paper* 14, 6 (2008), 1–36.

[9] Michal Herzenstein, Utpal M Dholakia, and Rick L Andrews. 2011. Strategic herding behavior in peer-to-peer loan auctions. *Journal of Interactive Marketing* 25, 1 (2011), 27–36.

[10] Michal Herzenstein, Scott Sonenshein, and Utpal M Dholakia. 2011. Tell me a good story and I may lend you money: The role of narratives in peer-to-peer lending decisions. *Journal of Marketing Research* 48, SPL (2011), S138–S149.

[11] Clayton J Hutto and Eric Gilbert. 2014. Vader: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth international AAAI conference on weblogs and social media*.

[12] Milad Malekipirbazari and Vural Aksakalli. 2015. Risk assessment in social lending via random forests. *Expert Systems with Applications* 42, 10 (2015), 4621–4631. https://doi.org/10.1016/j.eswa.2015.02.001

[13] John Neter, Michael H Kutner, Christopher J Nachtsheim, and William Wasserman. 1996. *Applied linear statistical models*. Vol. 4. Irwin Chicago.

[14] Prosper. 2018. Prosper Marketplace. https://www.prosper.com/invest. last accessed - 7/4/2019.

[15] Ke Ren and Avinash Malik. 2019. Investment Recommendation System for Low-Liquidity Online Peer to Peer Lending (P2PL) Marketplaces. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. ACM, 510–518.

[16] Joe Ryan, Katya Reuk, and Charles Wang. 2007. To fund or not to fund: Determinants of loan fundability in the prosper. com marketplace. *WP, The Standord Graduate School of Business* (2007).

[17] Joash Xu. 2015. Prosper Loan Data. https://github.com/joashxu/prosper-loan-data. last accessed - 15/9/2016.