# Lending Club

HarvardX - PH125.9x Data Science Capstone

*Emmanuel Rialland - https://github.com/Emmanuel_R8*

*November 03, 2019*

# Contents

# List of Tables

# List of Figures

# Introduction

Lending Club (*LC*) is an American company listed on the New York stock exchange that provides a platform for peer-to-peer lending. Unlike banks, it does not take deposits and invest them. It is purely a matching system. Each loan is split into $25 that multiple investors can invest in. LC is remunerated by fees received from both sides. LC states that they have intermediated more than $50bln since they started operations. Further description of the company is easily available online numerous sources.

In order for investors to make the best investment decisions, LC make a historical dataset publicly available to all registered investors. This dataset is the subject of this report. It was downloaded from the Kaggle data science website[1].

The size of the dataset is rich enough that it could be used to answer many different questions. We decided for a focused approach. Following Chapter 5 of (Peng, 2012), we will first formulate the question we want to answer to guide our analysis.

The business model of LC is to match borrowers and investors. Naturally, more people want to receive money than part with it. An important limiting factor to LC's growth is the ability to attract investors, build a trusting relationship where, as a minimum first step, investors trust LC to provide accurate, transparent and reliable information of the borrowers. For this purpose, LC decided not only to provide extensive information about potential borrowers' profile, but also historical information about past borrowers' performance. This is, as we understand, one of the key purposes of this dataset. We decided to use the dataset for this very purpose. Essentially, the questions are: **given a borrower profile, is his/her rating appropriate in terms of risk of default? And if a default occurs, what is the expected recovery? The summary question is: given a borrower profile, is the risk/reward balance appropriate to commit funds?** In answering this question, we understand that LC allows investment of very granular amounts. Therefore, even an individual investor can diversify his/her loan and risk portfolio. It is not necessary to 'gamble' funds on a single borrower. This is exactly what institutional investors achieve through syndication (although on a very different scale, typically $10-25mln for a medium-size bank).

For this exercise, we made two simplifying (hopefully not simplistic) assumptions:

- In determining the risk/return balance, we have not accounted for LC's cost of intermediation. By ignoring fees paid by both sides, we obviously overestimate the returns to the investors. But in first approximation, **we will assume that the risk/reward balance, from the investors' point of view, across ratings is independent from fees.** This is a simplification. Real-world fees are higher the lower the investment grade and push the investors to receive, and the borrowers to pay, higher interest margin.

- All-in interest rates paid by borrowers are fixed. This is highly desirable for borrowers to be able to manage their cashflow. However, an investor should always consider an investment return as a margin above a risk-free return. Banks would look at LIBOR; bond investors (e.g. life insurers) would look at government bonds. Those risk-free rates can change very quickly, whereas we understand that LC sets those rates on a

---

[1] https://www.kaggle.com/wendykan/lending-club-loan-data/data

less frequent basis. In other word, the risk premium will vary rapidly. **We assume that individual investors are 'in-elastic' to change in implied risk premia.** But we recognise this as a limitation of our work.

This report is organised as follows:

- [XXXX]

# Chapter 1

# Dataset

The data is sourced as a *SQLite* database that was imported as a `dataframe` with the `RSQLite` package. The variables were reformatted according to their respective types.

We also sourced US zip and FIPS codes, and macroeconomical data for possible geographical statistics. The source code for the data import and reformatting is given in appendix.

## 1.1 Preamble

The LendingClub dataset, although rich, is difficult to interpret. The only explanation of what the variables mean comes from a spreadsheet attached to the dataset. The explanations are not precise and/or subject to conflicting interpretation. Despite serching the LendingClub website, no further original information was found. We collected a number of reasonable assumptions in Appendix.

The dataset has been used a number of times in the past by various people. One paper (Kim and Cho, 2019) mentions they used a dataset that included 110 variables, which is less than ours with 145 variables. The dataset has changed over time in ways we do not know.

## 1.2 General presentation

The original dataset is rich: it includes 2260668 loan samples, each containing 144 variables (after the identification variables filled with null values). The loans were issued from 2007-06-01 to 2018-12-01.

### 1.2.1 Business volume

The dataset represents a total of ca.$34bln in loan principals, which is a substantial share of the total amount stated to have been intermediated to date by LC (reported to be $50bln+). About 60% of the portfolio is fully repaid. See Table 1.1.

Figure 1.5 plots the number, volume (cumulative principal amount) and average principal per loan. It shows that the business grew exponentially (in the common sense of the word) from inception until 2016. At this point, according to Wikipedia [1]:

" *Like other peer-to-peer lenders including Prosper, Sofi and Khutzpa.com, LendingClub experienced increasing difficulty attracting investors during early 2016. This led the firm to increase the interest rate it charges borrowers on three*

---

[1]source: https://en.wikipedia.org/wiki/LendingClub - Retrieval date 15 September 2019

Table 1.1: Number of loans per status

| Loan status | Count | Proportion (%) |
|---|---|---|
| Charged Off | 261655 | 11.574 |
| Current | 919695 | 40.682 |
| Default | 31 | 0.001 |
| Does not meet the credit policy. Status:Charged Off | 761 | 0.034 |
| Does not meet the credit policy. Status:Fully Paid | 1988 | 0.088 |
| Fully Paid | 1041952 | 46.090 |
| In Grace Period | 8952 | 0.396 |
| Late (16-30 days) | 3737 | 0.165 |
| Late (31-120 days) | 21897 | 0.969 |

*occasions during the first months of the year. The increase in interest rates and concerns over the impact of the slowing United States economy caused a large drop in LendingClub's share price."*

The number and volume of loans plotted have been aggregated by month. The growth is very smooth in the early years, and suddenly very volatile. As far as the first part of the dataset is concerned, a starting business could expect to be volatile and could witness a yearly cycle (expected from economic consumption figures) superimposed on the growth trend. This is not the case.

An interesting metric is that the average principal of loans has increased (see Figure **??**, on a sample of 100,000 loans). Partly, the increase in the early years could be interpreted success in confidence building. This metric plateau-ed in 2016 and decreased afterwards, but to a much lesser extent than the gross volume metrics. However, it is more volatile than the two previous metrics in the early years.

By the end of the dataset, all metrics have essentially recovered to their 2016 level.

### 1.2.2 Loan lifecyle and status

In the dataset, less loans are still outstanding than matured or "*charged off*" (term that LC use to mean partially or fully written off, i.e. there are no possibilty for LC and/or the investors to receive further payments). The share of outstanding loans is:

```
1  ## Share of current loans =  42.214 %
```

The dataset describes the life cycle of a loan. In the typical (ideal) case, we understand it to be:

Loan is approved → Full amount funded by investors → Loan marked as Current → Fully Paid

In the worst case, it is:

Loan is approved → Full amount funded by investors → Loan marked as Current →

→ Grace period (missed payments under 2 weeks) → Late 15 to 31 days →

→ Late 31 to 120 days → Default → Charged Off

Note that *Default* precedes and is distinct from *Charged Off* [2]. A couple of things could happen to a loan in default:

- LC and the borrower restructure the loan with a new repayment schedule, where the borrower may repay a lesser amount over a longer period; or,

- the claim could be sold to a debt recovery company that would buy the claim from LC/investors. This would be the final payment (if any) received by LC and the investors.

The dataset also describes situations where a borrower negotiated a restructuring of the repayment schedule in case of unexpected hardship (e.g. disaster, sudden unemployment).

Note that this progression of distinguishing default (event in time) and actual financial loss mirrors what banks and rating agencies do> The former is called the *Probability of Default* (PD), the latter *Loss Given Default* (LGD). Ratings change over time (in a process resembling a Markov Chains). LGD show some correlations with ratings. The dataset, although detailed, does not include the full life of each loan to conduct this sort of analysis (change of loan quality over time). This is an important reason why we decided to focus on the loan approval and expected return.

### 1.2.3   Loan application

Before a loan is approved, the borrower undergoes a review process that assess his/her capacity to repay. This includes:

- employment situation and income, as well whether this income and possibly its source has been independently verified;

- whether the application is made jointly (likely with a partner or a spouse, but there are no details);

- housing situation (owner, owner with current mortgage, rental) and in which county he/she lives (that piece of information is partially anonymised by removing the last 2 digits of the borrower's zipcode);

- the amount sought, its tenor and the purpose of the loan; and,

- what seems to be previous credit history (number of previous deliquencies). The dataset is very confusing in that regard: it is clear that such information relates to before the loan is approved in the case of the joint applicant. In the case of the principal borrower however, the variable descriptions could be read as pre-approval information, or information gathered during the life of the loan. We have assumed that the information related to the principal borrower is also pre-approval. We also used *Sales Supplements* from the LC website[3] that describe some of the information provided to investors. LendingClub also provides a summary description of its approval process in its regulatory filings with the Securities Exchange Commission (California, 2019).

### 1.2.4   Interest rates

Based on this information, the loan is approved or not. Approval includes the final amount (which could be lower than the amount requested), tenor (3 or 5 years) and a rating similar to those given to corporate borrowers. Unlike corporate borrowers however, the rating mechanically determines the rate of interest according to a grid known to the borrower in advance[4]. The rates have changed over time. Those changes where not as frequent as market conditions (e.g. changes in Federal Reserve Bank's rates)[5].

---

[2]See LendingClub FAQ at https://help.lendingclub.com/hc/en-us/articles/215488038

[3]See https://www.lendingclub.com/legal/prospectus

[4]https://www.lendingclub.com/investing/investor-education/interest-rates-and-fees

[5]Corporate borrowers would negociate interest margins on a case-by-case basis despite similar risk profiles.

Figure 1.1: Interest rates given rating

Figure 1.2: Interest rate per grade over time

Figure 1.1 [6] shows the predetermined interest rate depending on the initial rating as of July 2019.

At the date of this report, the ratings range from A (the best) down to D, each split in 5 sub-ratings. However, LC previously also intermediated loans rated F or G (until 6 November 2017) and E (until 30 June 2019) [7]. This explains that such ratings are in the dataset. We will assume that the ratings in the dataset are the rating at the time of approval and that, even if loans are re-rated by LC, the dataset does not reflect it.

Figures 1.2 shows the change in interest rate over time for different ratings and separated for each tenor. (Each figure is on a sample of 100,000 loans.) For each rating, we can see several parallel lines which correspond to the 5 sub-rating of each rating. We note that the range of interest rates has substantial widened over time. That is, the risk premium necessary to attract potential investors has had to substantially increase. In the most recent years, the highest rates exceed 30% which is higher than many credit cards.3-year loans are, unsurprisingly, considered safer (more A-rated, less G-rated). Identical ratings attract identical rates of interest.

By comparison, we plot the 3-year (in red) and 5-year (in blue) bank swap rates in Figure @(fig:swap-rates). We see that the swap curve has flattened in recent times (3-year and 5-y rates are almost identical). We also can see that in broad terms the interest rates charged reflect those underlying swap rates. It is therefore most relevant to examine the credit margins added to the swap rates.

Figures 1.4 shows the change in credit margin over time for different ratings and separated for each tenor. (Each figure is on a sample of 100,000 loans.) As above, for each rating, we can see several parallel lines which correspond to the 5 sub-rating of each rating. We note that the range of credit margins has widened over time but less than the interest rates. Identical ratings attract identical credit margins.

---

[6]source: https://www.lendingclub.com/investing/investor-education/interest-rates-and-fees
[7]See https://www.lendingclub.com/info/demand-and-credit-profile.action

Figure 1.3: Historical Swap Rates

Figure 1.4: Credit margins per grade over time

[TODO: DTI, amount... by grade]

### 1.2.5 Payments

The loans are approved for only two tenors, 3 and 5 years, with monthly repayments. Installments are calculated easily with the usual formula:

$$Installment = Principal \times \frac{1}{1 - \frac{1}{(1+rate)^N}}$$

Where $Principal$ is the amount borrowed, $rate = \frac{\text{Quoted Interest Rate}}{12}$ is the monthly interest rate, and $N$ is the number of installments (36 or 60 monthly payments). The following piece of code shows that the average error between this formula and the dataset value is about 2 cents. We therefore precisely understand this variable.

```
1   local({
2     installmentError <- loans %>%
3       mutate(
4         PMT = round(funded_amnt * int_rate / 12 / (1 - 1 / (1 + int_rate / 12) ^
5                                                 term), 2),
6         PMT_delta = abs(installment - PMT)
7       ) %>%
8       select(PMT_delta)
9
10    mean(100 * installmentError$PMT_delta)
11  })
```

## 1.3 Variables

We here present the dataset in a bit more details The full list of variable is given in appendix (see Table 2.1). This dataset will be reduced as we focused on our core question: *Are LC's loans priced appropriately?*.

### 1.3.1 General

### 1.3.2 Identification

The dataset is anonymised (all identifying ID numbers are deleted) and we therefore removed those columns from the dataset. Since the identification IDs have been removed to anonymise the dataset, we cannot see if a borrower borrowed several times.

Table 1.2: Matured loans per status

| Loan status | Count | Proportion (%) |
|---|---:|---:|
| Fully Paid | 1041952 | 26048800 |
| Charged Off | 261655 | 6541375 |
| Does not meet the credit policy. Status:Fully Paid | 1988 | 49700 |
| Does not meet the credit policy. Status:Charged Off | 761 | 19025 |

Figure 1.6: Funding and Write-offs by Sub-grades

## 1.4 Loan decision

As indicated in the introduction, our focus is on loans that have gone through their entire life cycle to consider their respective pricing, risk and profitability. To that effect, we will remove all loans which are still current (either performing or not). From here on, everything will be based on this reduced dataset.

In this reduced dataset, we focus on loans that have matured or been terminated. It contains 1306356 samples. Most of the loans (ca.80%) have been repaid in full. See Table 1.2.

When grouped by grade (Figure 1.6), we see a clear correlation between grade and default: the lower the grade the higher the portion defaults (all the way down to about 50%). In addition, most of the business is written in the B- or C-rating range.

# Chapter 2

# Modelling

At the outset, the dataset presents a number of challenges:

- There is a mix of continuous and categorical data.

- The number of observations is very large.

The diagram 2.1[1] shows a useful decision tree.

---

[1]Source: https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

Figure 2.1: Scikit Learn algorithm cheat-sheet

# Appendix

## 2.1 List of assumptions / limitations regarding the dataset

As mentioned during this report, we had to make numerous assumptions given the lack of clarity of the variable descriptions.

- The dataset does not contain any errors that we cannot notice (e.g. minor error of amount or rate, zipcode).

- The day-1 rating is between A1 and (and no lower than) G5. No note is rated lower than E5 from 6 November 2017, and lower than D5 from 30 June 2019.

- Credit history information for the principal borrower relates to pre-approval and not post-funding. This is clear for the joint applicants, but simply an assumption for the principal borrower.

- *Survival effect*: The dataset does not include applications that were rejected.

We do hope that LendingClub investors receive information of much better quality!

## 2.2 Data preparation and formatting

We used different sources of information:

- The LendingClub dataset made available on Kaggle;

- US georgraphical data about zip and FIPS codes;

- Market interest rates from the Saint Louis Federal Reserve Bank; and,

- Macro data from the same source.

We here show the code used to prepare the data.

### 2.2.1 LendinClub dataset

```
1   local({
2     #
3     # STEP 1: Download the dataset
4     #
5     #   Got to https://www.kaggle.com/wendykan/lending-club-loan-data
6     #
7     #   Download into the 'datasets' subdirectory
8     #   Unzip the file.
9     #   WARNING: The unzipping will be about 2.4GB
10    #
11    #   Name the sql database "datasets/lending_club.sqlite"
```

```
12    #
13
14    #
15    # STEP 2: Prepare the dabase as a tibble
16    #
17
18    library(RSQLite)
19    db_conn <-
20      dbConnect(RSQLite::SQLite(), "datasets/lending_club.sqlite")
21    dbListTables(db_conn)
22
23    # Returns a 2.96GB data frame
24    lending_club <- dbGetQuery(db_conn, "SELECT * FROM loan")
25    lending_club <- as_tibble(lending_club)
26
27    # Close the database
28    dbDisconnect(db_conn)
29
30    # Compressed to ca.285MB on disk
31    saveRDS(lending_club, "datasets/lending_club.rds")
32
33
34    library(tidyverse)
35    library(lubridate)
36    library(hablar)
37
38    # Before reformat in case the previous step was already done
39    # lending_club <- readRDS("datasets/lending_club.rds")
40    #
41    # str(lending_club)
42    #
43
44    # We leave the original dataset untouched and work with a copy.
45    lc <- lending_club
46
47    lc <- lc %>%
48      # Remove useless strings
49      mutate(
50        term       = str_remove(term, " months"),
51        emp_length = str_replace(emp_length, "<1", "0"),
52        emp_length = str_replace(emp_length, "10+", "10"),
53        emp_length = str_remove(emp_length, "years")
54      ) %>%
55
56      # Creates dates out of strings - Parse errors will be raised when no dates.
57      mutate(
58        debt_settlement_flag_date = as_date(dmy(
59          str_c("1-", debt_settlement_flag_date)
60        )),
61        earliest_cr_line          = as_date(dmy(str_c(
62          "1-", earliest_cr_line
63        ))),
64        hardship_start_date       = as_date(dmy(str_c(
65          "1-", hardship_start_date
66        ))),
```

```r
67        hardship_end_date        = as_date(dmy(str_c(
68          "1-", hardship_end_date
69        ))),
70        issue_d                  = as_date(dmy(str_c("1-", issue_d))),
71        last_credit_pull_d       = as_date(dmy(str_c(
72          "1-", last_credit_pull_d
73        ))),
74        last_pymnt_d             = as_date(dmy(str_c(
75          "1-", last_pymnt_d
76        ))),
77        next_pymnt_d             = as_date(dmy(str_c(
78          "1-", next_pymnt_d
79        ))),
80        payment_plan_start_date  = as_date(dmy(str_c(
81          "1-", payment_plan_start_date
82        ))),
83        sec_app_earliest_cr_line = as_date(dmy(str_c(
84          "1-", sec_app_earliest_cr_line
85        ))),
86        settlement_date          = as_date(dmy(str_c(
87          "1-", settlement_date
88        )))
89      ) %>%

91      # Bulk type conversion with convert from the `hablar` package
92      convert(
93        # Strings
94        chr(emp_title, title, url, zip_code),

96        # Factors
97        fct(
98          addr_state,
99          application_type,
100         debt_settlement_flag,
101         desc,
102         disbursement_method,
103         grade,
104         hardship_flag,
105         hardship_loan_status,
106         hardship_reason,
107         hardship_status,
108         hardship_type,
109         home_ownership,
110         id,
111         initial_list_status,
112         loan_status,
113         member_id,
114         policy_code,
115         purpose,
116         pymnt_plan,
117         settlement_status,
118         sub_grade,
119         verification_status,
120         verification_status_joint
121       ),
```

```
122
123        # Integers
124        int(
125          acc_now_delinq,
126          acc_open_past_24mths,
127          chargeoff_within_12_mths,
128          collections_12_mths_ex_med,
129          deferral_term,
130          delinq_2yrs,
131          emp_length,
132          hardship_dpd,
133          hardship_length,
134          inq_fi,
135          inq_last_12m,
136          inq_last_6mths,
137          mo_sin_old_il_acct,
138          mo_sin_old_rev_tl_op,
139          mo_sin_rcnt_rev_tl_op,
140          mo_sin_rcnt_tl,
141          mort_acc,
142          mths_since_last_delinq,
143          mths_since_last_major_derog,
144          mths_since_last_record,
145          mths_since_rcnt_il,
146          mths_since_recent_bc,
147          mths_since_recent_bc_dlq,
148          mths_since_recent_inq,
149          mths_since_recent_revol_delinq,
150          num_accts_ever_120_pd,
151          num_actv_bc_tl,
152          num_actv_rev_tl,
153          num_bc_sats,
154          num_bc_tl,
155          num_il_tl,
156          num_op_rev_tl,
157          num_rev_accts,
158          num_rev_tl_bal_gt_0,
159          num_sats,
160          num_tl_120dpd_2m,
161          num_tl_30dpd,
162          num_tl_90g_dpd_24m,
163          num_tl_op_past_12m,
164          open_acc,
165          open_acc_6m,
166          open_act_il,
167          open_il_12m,
168          open_il_24m,
169          open_rv_12m,
170          open_rv_24m,
171          sec_app_chargeoff_within_12_mths,
172          sec_app_collections_12_mths_ex_med,
173          sec_app_inq_last_6mths,
174          sec_app_mort_acc,
175          sec_app_mths_since_last_major_derog,
176          sec_app_num_rev_accts,
```

```
177          sec_app_open_acc,
178          sec_app_open_act_il,
179          term
180        ),
181
182        # Floating point
183        dbl(
184          all_util,
185          annual_inc,
186          annual_inc_joint,
187          avg_cur_bal,
188          bc_open_to_buy,
189          bc_util,
190          collection_recovery_fee,
191          delinq_amnt,
192          dti,
193          dti_joint,
194          funded_amnt,
195          funded_amnt_inv,
196          hardship_amount,
197          hardship_last_payment_amount,
198          hardship_payoff_balance_amount,
199          il_util,
200          installment,
201          int_rate,
202          last_pymnt_amnt,
203          loan_amnt,
204          max_bal_bc,
205          orig_projected_additional_accrued_interest,
206          out_prncp,
207          out_prncp_inv,
208          pct_tl_nvr_dlq,
209          percent_bc_gt_75,
210          pub_rec,
211          pub_rec_bankruptcies,
212          recoveries,
213          revol_bal,
214          revol_bal_joint,
215          revol_util,
216          sec_app_revol_util,
217          settlement_amount,
218          settlement_percentage,
219          tax_liens,
220          tot_coll_amt,
221          tot_cur_bal,
222          tot_hi_cred_lim,
223          total_acc,
224          total_bal_ex_mort,
225          total_bal_il,
226          total_bc_limit,
227          total_cu_tl,
228          total_il_high_credit_limit,
229          total_pymnt,
230          total_pymnt_inv,
231          total_rec_int,
```

```r
        total_rec_late_fee,
        total_rec_prncp,
        total_rev_hi_lim
      )
    ) %>%

    # Converts some values to 1/-1 (instead of Boolean)
    mutate(
      pymnt_plan =           if_else(pymnt_plan == "y",           1, -1),
      hardship_flag =        if_else(hardship_flag == "Y",        1, -1),
      debt_settlement_flag = if_else(debt_settlement_flag == "Y", 1, -1)
    ) %>%

    # Some values are percentages
    mutate(
      int_rate = int_rate / 100,
      dti = dti / 100,
      dti_joint = dti_joint / 100,
      revol_util = revol_util / 100,
      il_util = il_util / 100,
      all_util = all_util / 100,
      bc_open_to_buy = bc_util / 100,
      pct_tl_nvr_dlq = pct_tl_nvr_dlq / 100,
      percent_bc_gt_75 = percent_bc_gt_75 / 100,
      sec_app_revol_util = sec_app_revol_util / 100
    ) %>%

    # Create quasi-centered numerical grades out of grade factors with "A" = 3 down to "G" = -3
    mutate(grade_num = 4 - as.integer(grade)) %>%

    # Ditto with sub_grades. "A1" = +3.4, "A3" = +3.0, down to "G3" = -3.0, "G5" = -3.4
    mutate(sub_grade_num = 3.6 - as.integer(sub_grade) / 5) %>%

    # Keep the first 3 digits of the zipcode as numbers
    mutate(zip_code = as.integer(str_sub(zip_code, 1, 3))) %>%

    # Remove empty columns
    select(-id, -member_id, -url)

  saveRDS(lc, "datasets/lending_club_reformatted.rds")


  # Select loans which have matured or been terminated
  past_loans <- lc %>%
    filter(
      loan_status %in% c(
        "Charged Off",
        "Does not meet the credit policy. Status:Charged Off",
        "Does not meet the credit policy. Status:Fully Paid",
        "Fully Paid"
      )
    )

  saveRDS(past_loans, "datasets/lending_club_reformatted_paid.rds")
})
```

### 2.2.2 Zip codes and FIPS codes

The R package `zipcode` was installed.

```r
#
# ZIPCodes dataset.
#

library(zipcode)
data(zipcode)
zips <- zipcode %>%
  as_tibble() %>%
  mutate(zip = as.integer(str_sub(zip, 1, 3)))

saveRDS(zips, "datasets/zips.rds")
```

A csv file containing zip codes, FIPS codes and population information was downloaded from the *Simple Maps* [2] website.

```r
local({
  kaggleCodes <- read.csv("datasets/csv/ZIP-COUNTY-FIPS_2017-06.csv")

  kaggleCodes <-
    kaggleCodes %>%
    as_tibble() %>%
  mutate(zip = floor(ZIP/100),
         FIPS = STCOUNTYFP,
         COUNTYNAME = str_replace(COUNTYNAME, pattern = "County", replacement = ""),
         COUNTYNAME = str_replace(COUNTYNAME, pattern = "Borough", replacement = ""),
         COUNTYNAME = str_replace(COUNTYNAME, pattern = "Municipio", replacement = ""),
         COUNTYNAME = str_replace(COUNTYNAME, pattern = "Parish", replacement = ""),
         COUNTYNAME = str_replace(COUNTYNAME, pattern = "Census Area", replacement = "")) %>%
    rename(county = COUNTYNAME) %>%
    select(zip, county, FIPS) %>%
    arrange(zip)

  saveRDS(zipfips, "datasets/kaggleCodes.rds")
})
```

### 2.2.3 Market interest rates

Market interest rates (3-year and 5-year swap rates) were download from the Saint Louis Federal Reserve Bank. Datasets are split between before and after the LIBOR fixing scandal. The datasets are merged with disctinct dates.

Download sources are:

- Pre-LIBOR 3-y swap https://fred.stlouisfed.org/series/DSWP3

- Post-LIBOR 3-y swap https://fred.stlouisfed.org/series/ICERATES1100USD3Y

- Pre-LIBOR 5-y swap https://fred.stlouisfed.org/series/MSWP5

- Post-LIBOR 5-y swap https://fred.stlouisfed.org/series/ICERATES1100USD5Y

```r
local({
  LIBOR3Y <- read.csv("datasets/csv/DSWP3.csv") %>%
```

---

[2]https://simplemaps.com/data/us-zips

```r
    as_tibble() %>%
    filter(DSWP3 != ".") %>%
    mutate(DATE = as_date(DATE),
           RATE3Y = as.numeric(as.character(DSWP3)) / 100) %>%
    select(DATE, RATE3Y)

ICE3Y <- read.csv("datasets/csv/ICERATES1100USD3Y.csv") %>%
    as_tibble() %>%
    filter(ICERATES1100USD3Y != ".") %>%
    mutate(DATE = as_date(DATE),
           RATE3Y = as.numeric(as.character(ICERATES1100USD3Y)) / 100) %>%
    select(DATE, RATE3Y)


LIBOR5Y <- read.csv("datasets/csv/DSWP5.csv") %>%
    as_tibble() %>%
    filter(DSWP5 != ".") %>%
    mutate(DATE = as_date(DATE),
           RATE5Y = as.numeric(as.character(DSWP5)) / 100) %>%
    select(DATE, RATE5Y)

ICE5Y <- read.csv("datasets/csv/ICERATES1100USD5Y.csv") %>%
    as_tibble() %>%
    filter(ICERATES1100USD5Y != ".") %>%
    mutate(DATE = as_date(DATE),
           RATE5Y = as.numeric(as.character(ICERATES1100USD5Y)) / 100) %>%
    select(DATE, RATE5Y)

RATES3Y <- LIBOR3Y %>% rbind(ICE3Y) %>%
    arrange(DATE) %>% distinct(DATE, .keep_all = TRUE)

RATES5Y <- LIBOR5Y %>% rbind(ICE5Y) %>%
    arrange(DATE) %>% distinct(DATE, .keep_all = TRUE)

saveRDS(RATES3Y, "datasets/rates3Y.rds")
saveRDS(RATES5Y, "datasets/rates5Y.rds")


# Note there are 7212 days from 1 Jan 2000 to 30 Sep 2019
#
# (ymd("2000-01-01") %--% ymd("2019-09-30")) %/% days(1)
RATES <- tibble(n = seq(0, 7212)) %>%

    # Create a column with all dates
    mutate(DATE = ymd("2000-01-01") + days(n)) %>%
    select(-n) %>%

    # Add all daily 3- then 5-year rates and fill missing down
    left_join(RATES3Y) %>%
    fill(RATE3Y, .direction = "down") %>%

    left_join(RATES5Y) %>%
    fill(RATE5Y, .direction = "down")

saveRDS(RATES, "datasets/rates.rds")
```

```
58      })
```

### 2.2.4   Macro-economical data

Macro-economical datasets were sourced from the same website as Microsoft Excel files. They were converted as-is to tab-separated csv files with LibreOffice.

- Median income per household: https://geofred.stlouisfed.org/map/?th=pubugn&cc=5&rc=false&im=
  fractile&sb&lng=-112.41&lat=44.31&zm=4&sl&sv&am=Average&at=Not%20Seasonally%20Adjusted,
  %20Annual,%20Dollars&sti=2022&fq=Annual&rt=county&un=lin&dt=2017-01-01

- Per capita personal income: https://geofred.stlouisfed.org/map/?th=pubugn&cc=5&rc=false&im=
  fractile&sb&lng=-112.41&lat=44.31&zm=4&sl&sv&am=Average&at=Not%20Seasonally%20Adjusted,
  %20Annual,%20Dollars&sti=882&fq=Annual&rt=county&un=lin&dt=2017-01-01

- Unemployment: https://geofred.stlouisfed.org/map/?th=rdpu&cc=5&rc=false&im=fractile&sb&lng=
  -90&lat=40&zm=4&sl&sv&am=Average&at=Not%20Seasonally%20Adjusted,%20Monthly,%20Percent&
  sti=1224&fq=Monthly&rt=county&un=lin&dt=2019-08-01

```r
1   local({
2     ###############################################################################################
3     ##
4     ## Median income per household by FIPS from 2002 to 2017
5     ##
6     # Prepare median income
7     medianIncome <-
8       # Load the dataset after dropping the first line
9       read.csv(
10        "datasets/csv/GeoFRED_Estimate_of_Median_Household_Income_by_County_Dollars.csv",
11        sep = "\t",
12        skip = 1,
13        stringsAsFactors = FALSE
14      ) %>%
15
16      # Drops columnsn containing a unique identifier and the FIPS name
17      select(-"Series.ID", -"Region.Name") %>%
18
19      # Rename the relevant column to 'FIPS'
20      rename(FIPS = "Region.Code") %>%
21
22      # Order by FIPS
23      arrange(FIPS) %>%
24
25      # Convert to a 'long' table, i.e. one column for FIPS, one for date, one for income
26      pivot_longer(cols = starts_with("X"),
27                   names_to = "Date",
28                   values_to = "medianIncome") %>%
29
30      # Create actual dates
31      mutate(Date = str_replace(Date, "[X]", ""),
32             Date = ymd(str_c(Date, "-12-31")))
33
34      saveRDS(medianIncome, "datasets/medianincome.rds")
35
36
37
```

```
38  ################################################################################
39  ##
40  ## Per capita income by FIPS from 2002 to 2017
41  ##
42  personalIncome <-
43    # Load the dataset after dropping the first line
44    read.csv(
45      "datasets/csv/GeoFRED_Per_Capita_Personal_Income_by_County_Dollars.csv",
46      sep = "\t",
47      skip = 1,
48      stringsAsFactors = FALSE
49    ) %>%
50
51    # Drops columnsn containing a unique identifier and the FIPS name
52    select(-"Series.ID", -"Region.Name") %>%
53
54    # Rename the relevant column to 'FIPS'
55    rename(FIPS = "Region.Code") %>%
56
57    # Order by FIPS
58    arrange(FIPS) %>%
59
60    # Convert to a 'long' table, i.e. one column for FIPS, one for date, one for income
61    pivot_longer(cols = starts_with("X"),
62                 names_to = "Date",
63                 values_to = "personalIncome") %>%
64
65    # Create actual dates
66    mutate(Date = str_replace(Date, "[X]", ""),
67           Date = ymd(str_c(Date, "-12-31")))
68
69  saveRDS(personalIncome, "datasets/personalincome.rds")
70
71
72  ################################################################################
73  ##
74  ## Unemplyment rate monthly by FIPS from January 2000 to August 2019
75  ##
76  unemploymentRate <-
77    # Load the dataset after dropping the first line
78    read.csv(
79      "datasets/csv/GeoFRED_Unemployment_Rate_by_County_Percent.csv",
80      sep = "\t",
81      skip = 1,
82      stringsAsFactors = FALSE
83    ) %>%
84
85    # Drops columnsn containing a unique identifier and the FIPS name
86    select(-"Series.ID", -"Region.Name") %>%
87
88    # Rename the relevant column to 'FIPS'
89    rename(FIPS = "Region.Code") %>%
90
91    # Order by FIPS
92    arrange(FIPS) %>%
```

```
93
94      mutate_all(as.double) %>%
95
96      # Convert to a 'long' table, i.e. one column for FIPS, one for date, one for income
97      pivot_longer(cols = starts_with("X"),
98                   names_to = "Date",
99                   values_to = "unemploymentRate",
100                  values_ptypes = c("unemploymentRate", numeric)) %>%
101
102     # Converts the content of the Year column to an actual date
103     mutate(
104       Date = str_replace(Date, "[X]", ""),
105       Date = str_replace(Date, "[.]", "-"),
106       Date = ymd(str_c(Date, "-1"))
107     )
108
109   saveRDS(unemploymentRate, "datasets/unemployment.rds")
110 })
```

## 2.3    List of variables

This table presents the list of variables provided in the original dataset. The descriptions come from a spreadsheet attached with the dataset and, unfortunately, are not extremely precise and subject to interpretation. We added comments and/or particular interpretations in *CAPITAL LETTERS*.

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle

| Variable Name | Description |
| --- | --- |
| loan_amnt | The listed amount of the loan applied for by the borrower. If at some point in time, the credit department reduces the loan amount, then it will be reflected in this value. |
| funded_amnt | The total amount committed to that loan at that point in time. |
| funded_amnt_inv | The total amount committed by investors for that loan at that point in time. |
| term | The number of payments on the loan. Values are in months and can be either 36 or 60. |
| int_rate | Interest Rate on the loan |
| installment | The monthly payment owed by the borrower if the loan originates. |
| grade | LC assigned loan grade |
| sub_grade | LC assigned loan subgrade |
| emp_title | The job title supplied by the Borrower when applying for the loan. |
| emp_length | Employment length in years. Possible values are between 0 and 10 where 0 means less than one year and 10 means ten or more years. |
| home_ownership | The home ownership status provided by the borrower during registration or obtained from the credit report. Our values are: RENT, OWN, MORTGAGE, OTHER, NONE |

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle *(continued)*

| Variable Name | Description |
| --- | --- |
| annual_inc | The self-reported annual income provided by the borrower during registration. NOT USED AS A VARIABLE SINCE JOINT INCOME ALREADY INCLUDES IT. |
| verification_status | Indicates if income was verified by LC, not verified, or if the income source was verified |
| issue_d | The month which the loan was funded |
| loan_status | Current status of the loan |
| pymnt_plan | Indicates if a payment plan has been put in place for the loan |
| url | URL for the LC page with listing data. |
| desc | Loan description provided by the borrower |
| purpose | A category provided by the borrower for the loan request. |
| title | The loan title provided by the borrower |
| zip_code | The first 3 numbers of the zip code provided by the borrower in the loan application. |
| addr_state | The state provided by the borrower in the loan application |
| dti | A ratio calculated using the borrower s total monthly debt payments on the total debt obligations, excluding mortgage and the requested LC loan, divided by the borrower s self-reported monthly income. NOT USED AS A VARIABLE. ONLY USE JOINT DTI. |
| delinq_2yrs | The number of 30+ days past-due incidences of delinquency in the borrower s credit file for the past 2 years |
| earliest_cr_line | The month the borrower s earliest reported credit line was opened |
| inq_last_6mths | The number of inquiries in past 6 months (excluding auto and mortgage inquiries) |
| mths_since_last_delinq | The number of months since the borrower s last delinquency. |
| mths_since_last_record | The number of months since the last public record. |
| open_acc | The number of open credit lines in the borrower s credit file. |
| pub_rec | Number of derogatory public records |
| revol_bal | Total credit revolving balance |
| revol_util | Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit. |
| total_acc | The total number of credit lines currently in the borrower s credit file |
| initial_list_status | The initial listing status of the loan. Possible values are – W, F |
| out_prncp | Remaining outstanding principal for total amount funded |
| out_prncp_inv | Remaining outstanding principal for portion of total amount funded by investors |

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle *(continued)*

| Variable Name | Description |
| --- | --- |
| total_pymnt | Payments received to date for total amount funded |
| total_pymnt_inv | Payments received to date for portion of total amount funded by investors |
| total_rec_prncp | Principal received to date |
| total_rec_int | Interest received to date |
| total_rec_late_fee | Late fees received to date |
| recoveries | Post charge off gross recovery |
| collection_recovery_fee | Post charge off collection fee |
| last_pymnt_d | Last month payment was received |
| last_pymnt_amnt | Last total payment amount received |
| next_pymnt_d | Next scheduled payment date |
| last_credit_pull_d | The most recent month LC pulled credit for this loan |
| collections_12_mths_ex_med | Number of collections in 12 months excluding medical collections |
| mths_since_last_major_derog | Months since most recent 90-day or worse rating |
| policy_code | Publicly available policy_code=1 / New products not publicly available policy_code=2 |
| application_type | Indicates whether the loan is an individual application or a joint application with two coborrowers |
| annual_inc_joint | The combined self-reported annual income provided by the coborrowers during registration |
| dti_joint | A ratio calculated using the coborrowers total monthly payments on the total debt obligations, excluding mortgages and the requested LC loan, divided by the coborrowers combined self-reported monthly income |
| verification_status_joint | Indicates if income was verified by LC, not verified, or if the income source was verified |
| acc_now_delinq | The number of accounts on which the borrower is now delinquent. |
| tot_coll_amt | Total collection amounts ever owed |
| tot_cur_bal | Total current balance of all accounts |
| open_acc_6m | Number of open trades in last 6 months |
| open_act_il | Number of currently active installment trades |
| open_il_12m | Number of installment accounts opened in past 12 months |
| open_il_24m | Number of installment accounts opened in past 24 months |
| mths_since_rcnt_il | Months since most recent instalment accounts opened |
| total_bal_il | Total current balance of all installment accounts |
| il_util | Ratio of total current balance to high credit/credit limit on all install acct |
| open_rv_12m | Number of revolving trades opened in past 12 months |
| open_rv_24m | Number of revolving trades opened in past 24 months |
| max_bal_bc | Maximum current balance owed on all revolving accounts |

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle *(continued)*

| Variable Name | Description |
| --- | --- |
| all_util | Balance to credit limit on all trades |
| total_rev_hi_lim | Total revolving high credit/credit limit |
| inq_fi | Number of personal finance inquiries |
| total_cu_tl | Number of finance trades |
| inq_last_12m | Number of credit inquiries in past 12 months |
| acc_open_past_24mths | Number of trades opened in past 24 months. |
| avg_cur_bal | Average current balance of all accounts |
| bc_open_to_buy | Total open to buy on revolving bankcards. |
| bc_util | Ratio of total current balance to high credit/credit limit for all bankcard accounts. |
| chargeoff_within_12_mths | Number of charge-offs within 12 months |
| delinq_amnt | The past-due amount owed for the accounts on which the borrower is now delinquent. |
| mo_sin_old_il_acct | Months since oldest bank instalment account opened |
| mo_sin_old_rev_tl_op | Months since oldest revolving account opened |
| mo_sin_rcnt_rev_tl_op | Months since most recent revolving account opened |
| mo_sin_rcnt_tl | Months since most recent account opened |
| mort_acc | Number of mortgage accounts. |
| mths_since_recent_bc | Months since most recent bankcard account opened. |
| mths_since_recent_bc_dlq | Months since most recent bankcard delinquency |
| mths_since_recent_inq | Months since most recent inquiry. |
| mths_since_recent_revol_delinq | Months since most recent revolving delinquency. |
| num_accts_ever_120_pd | Number of accounts ever 120 or more days past due |
| num_actv_bc_tl | Number of currently active bankcard accounts |
| num_actv_rev_tl | Number of currently active revolving trades |
| num_bc_sats | Number of satisfactory bankcard accounts |
| num_bc_tl | Number of bankcard accounts |
| num_il_tl | Number of installment accounts |
| num_op_rev_tl | Number of open revolving accounts |
| num_rev_accts | Number of revolving accounts |
| num_rev_tl_bal_gt_0 | Number of revolving trades with balance >0 |
| num_sats | Number of satisfactory accounts |
| num_tl_120dpd_2m | Number of accounts currently 120 days past due (updated in past 2 months) |
| num_tl_30dpd | Number of accounts currently 30 days past due (updated in past 2 months) |
| num_tl_90g_dpd_24m | Number of accounts 90 or more days past due in last 24 months |
| num_tl_op_past_12m | Number of accounts opened in past 12 months |
| pct_tl_nvr_dlq | Percent of trades never delinquent |
| percent_bc_gt_75 | Percentage of all bankcard accounts > 75% of limit. |
| pub_rec_bankruptcies | Number of public record bankruptcies |
| tax_liens | Number of tax liens |
| tot_hi_cred_lim | Total high credit/credit limit |
| total_bal_ex_mort | Total credit balance excluding mortgage |
| total_bc_limit | Total bankcard high credit/credit limit |
| total_il_high_credit_limit | Total installment high credit/credit limit |
| revol_bal_joint | Total credit revolving balance |

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle *(continued)*

| Variable Name | Description |
|---|---|
| sec_app_earliest_cr_line | Earliest credit line at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_inq_last_6mths | Credit inquiries in the last 6 months at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_mort_acc | Number of mortgage accounts at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_open_acc | Number of open trades at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_revol_util | Ratio of total current balance to high credit/credit limit for all revolving accounts. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_open_act_il | Number of currently active installment trades at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_num_rev_accts | Number of revolving accounts at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_chargeoff_within_12_mths | Number of charge-offs within last 12 months at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_collections_12_mths_ex_med | Number of collections within last 12 months excluding medical collections at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| sec_app_mths_since_last_major_derog | Months since most recent 90-day or worse rating at time of application for the secondary applicant. VARIABLE NOT USED. WE RELY ON THE MAIN BORROWER IN THE FIRST INSTANCE. |
| hardship_flag | Flags whether or not the borrower is on a hardship plan |
| hardship_type | Describes the hardship plan offering |
| hardship_reason | Describes the reason the hardship plan was offered |
| hardship_status | Describes if the hardship plan is active, pending, cancelled, completed, or broken |
| deferral_term | Amount of months that the borrower is expected to pay less than the contractual monthly payment amount due to a hardship plan |

Table 2.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle *(continued)*

| Variable Name | Description |
| --- | --- |
| hardship_amount | The interest payment that the borrower has committed to make each month while they are on a hardship plan |
| hardship_start_date | The start date of the hardship plan period |
| hardship_end_date | The end date of the hardship plan period |
| payment_plan_start_date | The day the first hardship plan payment is due. For example, if a borrower has a hardship plan period of 3 months, the start date is the start of the three-month period in which the borrower is allowed to make interest-only payments. |
| hardship_length | The number of months the borrower will make smaller payments than normally obligated due to a hardship plan |
| hardship_dpd | Account days past due as of the hardship plan start date |
| hardship_loan_status | Loan Status as of the hardship plan start date |
| orig_projected_additional_accrued_interest | The original projected additional interest amount that will accrue for the given hardship payment plan as of the Hardship Start Date. This field will be null if the borrower has broken their hardship payment plan. |
| hardship_payoff_balance_amount | The payoff balance amount as of the hardship plan start date |
| hardship_last_payment_amount | The last payment amount as of the hardship plan start date |
| disbursement_method | The method by which the borrower receives their loan. Possible values are: CASH, DIRECT_PAY |
| debt_settlement_flag | Flags whether or not the borrower, who has charged-off, is working with a debt-settlement company. |
| debt_settlement_flag_date | The most recent date that the Debt_Settlement_Flag has been set |
| settlement_status | The status of the borrower's settlement plan. Possible values are: COMPLETE, ACTIVE, BROKEN, CANCELLED, DENIED, DRAFT |
| settlement_date | The date that the borrower agrees to the settlement plan |
| settlement_amount | The loan amount that the borrower has agreed to settle for |
| settlement_percentage | The settlement amount as a percentage of the payoff balance amount on the loan |
| settlement_term | The number of months that the borrower will be on the settlement plan |

## 2.4 System version

```
1  ##                          sysname
2  ##                          "Linux"
3  ##                          release
4  ##              "5.3.0-21-generic"
```

```
##                          version
## "#22-Ubuntu SMP Tue Oct 29 22:55:51 UTC 2019"
##                         nodename
##                          "x260"
##                          machine
##                         "x86_64"
##                           login
##                        "unknown"
##                            user
##                       "emmanuel"
##                   effective_user
##                       "emmanuel"
```

# Bibliography

California, L. C. S. F. (2019). Prospectus Regulatory Filing S3-ASR for Member Payment Dependent Notes. https://www.sec.gov/Archives/edgar/data/1409970/000140997019000988/0001409970-19-000988-index.htm. [Note: Accessed 31 October 2019].

Kim, A. and Cho, S.-B. (2019). An ensemble semi-supervised learning method for predicting defaults in social lending. *Engineering Applications of Artificial Intelligence*, 81:193–199.

Peng, R. (2012). *Exploratory data analysis with R*. Lulu. com.