

LendingClub Loans Pricing

HarvardX - PH125.9x Data Science Capstone

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

December 11, 2019

Contents

Introduction	5
1 Internal Rate of Return, Credit Margins and Net Present Values	7
1.1 Background	7
1.2 Internal Rate of Return	8
1.3 Dataset calculation	9
1.3.1 IRR	9
1.3.2 Credit Margins	9
1.3.3 NPV	9
2 Dataset	10
2.1 Preamble	10
2.2 General presentation	10
2.2.1 Business volume	10
2.2.2 Loan lifecycle and status	12
2.2.3 Loan application	13
2.3 Rates	13
2.3.1 IRR and required credit margins	13
2.3.2 Choice of predictors	14
2.3.3 Interest rates	14
2.3.4 Purpose	17
2.3.5 Payments	17
2.4 Net present value	19
2.4.1 Average NPV and credit margin by subgrade	19
2.4.2 Distribution of principal losses by rating	19
2.4.3 NPV distribution by rating	19
2.5 Loan decision	21
3 Logistic Regression Model and Credit Scorecard	25
3.1 Logistic Regression	26
3.1.1 Data preparation	26
3.2 Binning and Weight of Evidence	28
3.2.1 Background	29
3.2.2 Binning	29
3.2.3 WOE and IV definitions	29
3.2.4 Calculating WoE and IV	31
3.2.5 Loop through all variables	32
3.2.6 Select relevant variables	33
3.2.7 Create data table with one-hot encoding	36
3.2.8 Comparison of individual characteristics	38
3.3 Logistic Regression	43

3.3.1	Remove identical bins	43
3.3.2	First model - complete dataset	44
3.3.3	Second training - fitted model less colinear bins	46
3.3.4	Third training - second fitted model less non-significant bins	48
3.4	Model results	50
3.4.1	Final list of variables	50
3.4.2	Scoring	51
3.4.3	Model scorecard	52
3.5	Training set	54
3.5.1	Densities of the training results	55
3.6	Test set	56
3.7	Confusion matrix	59
3.8	ROC Curve	61
4	Conclusion	65
5	Errands and post-mortem	67
5.1	Modeling	67
5.2	Geographical data	67
5.3	Stochastic Gradient Descent	68
5.3.1	Conclusions	68
5.3.2	Gradient descent	69
5.3.3	Stochastic Gradient Descent	70
5.3.4	Cost function for multi-modal NPV	70
5.4	Final final Conclusion	71
6	Appendix	72
6.1	List of assumptions / limitations regarding the dataset	72
6.2	Data preparation and formatting	72
6.2.1	LendingClub dataset	73
6.2.2	Zip codes and FIPS codes	73
6.2.3	Market interest rates	73
6.3	List of variables	73
6.4	Maxima derivation of the cost function	80
6.5	System version	81

List of Tables

2.1	Number of loans per status	11
2.2	Number of loans per purpose	18
2.3	Matured loans per status	22
3.1	Variable relevance by Information Value	34
3.2	15 top variables by IV	34
3.3	15 used top informative variables by IV	35
3.4	15 least informative variables by IV	36
3.5	First model. Example of NA coefficients.	45
3.6	First training. 15 best bins	46
3.7	Second training. 15 best bins	47
3.8	Third training: 15 best bins	49
3.9	Training AIC and Log-likelihood	50
6.1	Description of the dataset variables as provided in the dataset downloaded from Kaggle	74

List of Figures

2.1	Business volume written per month	11
2.2	Credit margins per grade over time	13
2.3	Credit margins per grade over time	14
2.4	Credit margins per grade over time	15
2.5	Interest rates given rating	15
2.6	Interest rate per grade over time	16
2.7	Historical Swap Rates	16
2.8	Credit margins per grade over time	17
2.9	Histograms of credit margins per purpose	17
2.10	Boxplots of credit margins per purpose	18
2.11	Average NPV et credit margin (%) depending on sub-rating	20
2.12	Distribution of the Principal Loss (%) depending on rating (y-axis square-root scaling)	20
2.13	Distribution of NPV (%) depending on rating (y-axis square-root scaling)	21
2.14	NPV % higher than 120% (no y-axis scaling)	22
2.15	NPV % around 41% (no y-axis scaling)	22
2.16	NPV % around -1% (no y-axis scaling)	23
2.17	NPV % for close to total loss % (no y-axis scaling)	23
2.18	Funding and Write-offs by Sub-grades	24
3.1	Model results on the training set	56
3.2	Density of loans by scorecard	58
3.3	Logit value predicted by the model	58
3.4	Probability of a loan being GOOD predicted by the model	59
3.5	Odds of a loan being GOOD predicted by the model	60
3.6	Receiver Operating Characteristic	62
3.7	Precision/Recall curve	62
3.8	Sensitivity/Specificity curve	63
3.9	Lift Chart	64
5.1	Scikit Learn algorithm cheat-sheet	69

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 state they have intermediated more than \$50bln since they started operations. Further description of the company is easily available from numerous online sources.

The business model of LC is to match borrowers and investors. Naturally, people prefer receiving money to parting 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 publicly available to all registered investors. This dataset is the subject of this report. It was downloaded from the Kaggle data science website¹.

The size of the dataset is rich enough that it could be used to answer many different questions. The questions could be: **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 the course of preparing this report, we considered those and eventually focused on researching the first question lenders seek to answer: **What is the probability of a borrower no repaying principal with interest in full?** Following Chapter 5 of (Peng, 2012), formulating this question will guide our analysis.

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). An investor can diversify his/her portfolio of loans across many borrowers so that probabilities of default can be considered on a statistical basis.

This report is organised as follows:

- We first introduce some financial terms and concepts that will be used in the rest of the report. This can be skipped if you studied finance. Those are just basic concepts.
- The second section introduces the dataset and uses a number of visualisations to illustrate some important aspects. We also provide the calculation of some financial amounts introduced in the first section.
- The third section described the model we used to assess the probability of a loan defaulting based on the information provided by an applicant. This is just one example of the sort of

¹<https://www.kaggle.com/wendykan/lending-club-loan-data/data>

questions one could try to answer based on the dataset.

- Finally, we assess our results in the conclusion. We provided numerous avenues to improve on this project.
- A post-mortem section is provided for interest purposes only.

IMPORTANT WARNING:

The calculations presented here are simplistic, although they bear some resemblance to what financial institutions (FIs) do. The literature on credit assessment and pricing is very rich and very complex. Finding the optimal capital allocation to particular risks while at the same time satisfying internal risk policies and regulatory requirements is a problem that financial institutions have yet to solve in full. Investing in a loan is not only a matter of assessing the risk of a particular borrower, but also assessing systemic risks (which exist across all borrowers), risks associated with funding the loan (interest, currency and liquidity markets), each requiring a risk assessment and pricing.

In other words, nobody would, let alone should, make any investment decision based on the calculations below.

Chapter 1

Internal Rate of Return, Credit Margins and Net Present Values

In this section, we introduce basic financial concepts. We also describe variables that we will generate for each loan and that will be used in the rest of this report.

1.1 Background

This subsection can be skipped by anybody with basic financial knowledge.

A bird in hand or two in the bush; a penny today or a pound tomorrow. What is the price of delaying obtaining and owning something? This is what pricing a loan is about. A lender could keep his/her cash in hand, or lend it and have it in hand later. He/she would accept this in exchange for receiving a bit more: this is the rate of interest. A lender wants to be compensated for delaying the possibility of using the cash, but also for taking the risk of not receiving it, partially or in full, when repayment is due.

There are borrowers that one can see as (almost) completely safe or risk-free such as central banks or governments of strong economies. A lender always has the possibility to lend to them instead of more risky borrowers. Therefore, a lender would require a higher interest rate than risk-free. The additional interest that a lender requires is commensurate with the risk of the borrower not repaying (*credit worthiness*) and is called the *credit margin*.

For each individual borrower, an FI would assess information provided by the borrower and historical data to answer the question: considering historical borrowers with a profile similar to the applicant's, what is the probability of not getting principal and interest back (*Probability of Default* or *PD*)? And, in case the borrower stops paying and, maybe using additional courses of action (such as seizing and selling assets), what is the total loss that could be expected on average (*Loss given Default* or *LGD*)?

Making that assessment, the FI would require an interest rate which would roughly be the sum of:

- the risk-free rate;
- a margin to cover the average loss of similar individual borrowers¹;

¹It is important to realise that the average margin only brings the borrower back to having earned the risk-free rate average: the additional income from the credit margin will be spent to cover average losses. In addition, we present the credit margin as income against borrower-specific losses. It does not address a lot of other risks such as correlation risks: a borrower might default because the economy as a whole gets worse, in which case many borrowers will default. This is a cyclical *systemic* risk similar to the 2007 US real estate crisis.

- a margin to cover all the operational costs of running their operations; and,
- a margin to remunerate the capital allocated by the FI (banking regulations require all banks to allocate an amount of capital against any risk taken; this is stipulated in a number of complex rules).

Said crudely, this total is the amount for the FI to get out of bed and look at a loan. Although this sounds like an exact science (for some definition of the word), it is not. At the end of the day, the FI will also have to contend with the competition from other FIs or non banking lenders, market liquidity (if there is a lot of money available to be lent, it brings prices down) and, of course, whether the borrower would at all be interested in accepting that cost.

Note that the dataset is distorted by this additional survival effect: the application information of many loans does not appear merely because the rate of interest was considered too high (this is not dissimilar to *survival effects* where some data did not survive through the history of a dataset²).

1.2 Internal Rate of Return

For the purpose of this report, we will simplify things: we will only consider the first two components of the interest rate. The risk-free rate and the credit margin that would cover the cost of default/losses of individual borrowers.

With respect to a given loan and its cash flow, two calculations are important: the Net Present Value (*NPV*) and the Internal Rate of Return (*IRR*). If we remember that an FI is indifferent to holding a principal P today or receiving it a year later with an annual interest (i.e. $P \times (1 + r)$ where r is the annual rate of interest), we can say that any amount CF_1 received in a year is equivalent to $CF_0 = \frac{CF_1}{1+r}$ today. More generally, a stream of future cash receipts is worth:

$$NPV(r) = \sum_{Year=1}^{Year=n} \frac{CF_i}{(1+r)^i}$$

The amount $NPV(r)$ is called the *Net Present Value* of the cash flow discounted at the rate r . Given that the LendingClub repayments are monthly, the formula becomes:

$$NPV(r) = \sum_{Month_i=1}^{Month_i=12 \times n} \frac{CF_i}{(1 + \frac{r}{12})^i}$$

If we now have a day 1 cash flow CF_0 , we can calculate:

$$CF_0 - \sum_{Year=1}^{Year=n} \frac{CF_i}{(1+r)^i}$$

However, for any given CF_0 , there is no reason that it would equal the NPV of the future cash flow (i.e no reason why the difference would be equal to zero). But this equation is a function of r . If we can find a value of r that zeroes this formula, it is called the internal rate of return of the cash flow³:

²A well-known example is historical stock prices which disappear when companies are de-listed or go bankrupt.

³Given the coefficients of such polynomials, we never get strange root values.

$$CF_0 - \sum_{Year=1}^{Year=n} \frac{CF_i}{(1 + IRR(CF))^i} = 0$$

or for monthly cash flow:

$$CF_0 - \sum_{Month=1}^{Month=12 \times n} \frac{CF_i}{(1 + \frac{IRR(CF)}{12})^i} = 0$$

1.3 Dataset calculation

For each loan, we calculated the IRR, credit margin and NPV. The calculations were performed in Julia due to R's slow performance on such a large dataset ⁴.

The source code files are in the `Scripts` folder.

1.3.1 IRR

We used the dataset to calculate the IRR of each loan. We used the following information for the dataset: `funded_amnt` (loan amount funded), `int_rate` (all-in interest rate), `term` (tenor of the loan in months), `total_pymnt` (total cumulative amount received from the borrower), `total_rec_prncp` (amount repaid allocated to principal repayment), `total_rec_int` (amount repaid allocated to interest payment), `recoveries` (any amount recovered later from the borrower) and `total_rec_late_fee` (any late payments fees paid by the borrower).

From that information, we recreated a cash flow for each loan. Unfortunately, this R code takes close to a full day to run on the entire dataset of completed loans (i.e. excluding all ongoing loans). This is just impractical for anybody to run to check this report and the resulting IRR results dataset is included in the Github repository. To make things practical, the dataset was actually created using code in Julia⁵. It is a direct translation of the R code, with a similar syntax (therefore very easy to follow). The Julia code runs about 500 times quicker (this is not a typo), or about 3 minutes. We appreciate that this is the departure from the assignment description.

Similarly, we calculate the credit margin required by each loan noting that Risk-free + Credit Margin = $IRR(loan)$.

1.3.2 Credit Margins

As noted in the previous section, risk-free rates change over time. When solving for the credit margin, we use the relevant risk-free rate.

Again, this was coded in Julia. The Julia code here takes about 1h20min to run. On the assumptions that the equivalent R code would take therefore almost 30 days to run through the full dataset, we did not write any R code for this calculation.

1.3.3 NPV

We also calculated the NPV of each loan, both as an absolute dollar amount and as a portion of the original. Again, this was coded in Julia. Visualisations based on the credit margins and NPV are in the next section.

⁴The resulting datasets are too large to be made available on the GitHub repository.

⁵<https://julialang.org/>

Chapter 2

Dataset

The data is sourced as a *SQLite* database that downloaded from teh Kaggle website (Kan, 2019) and imported as a `tibble` `dataframe` with the `RSQLite` package. The variables were reformatted according to their respective types. The full list of variable is given in Appendix (see Table ??). This dataset will be reduced as we focused on the core intent of modeling the probability of default.

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

2.1 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 searching the LendingClub website, no further original information was found. We collected a number of reasonable assumptions in Appendix (see subsection 6.1 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. It is therefore clear that the dataset has changed over time in ways we do not know. For example, have loans been excluded because the full 145 variables were not available?

2.2 General presentation

The original dataset is large: it includes 2260668 loan samples, each containing 145 variables (after the identification variables filled with null values). The loans were issued from 2007-06-01 to 2018-12-01.

2.2.1 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 (publicly reported to be \$50bln+). About 55%/60% of the portfolio is not current anymore (either fully or partially repaid). See Table 2.1.

Figure 2.1 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

Table 2.1: 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

inception until 2016. At this point, according to Wikipedia ¹:

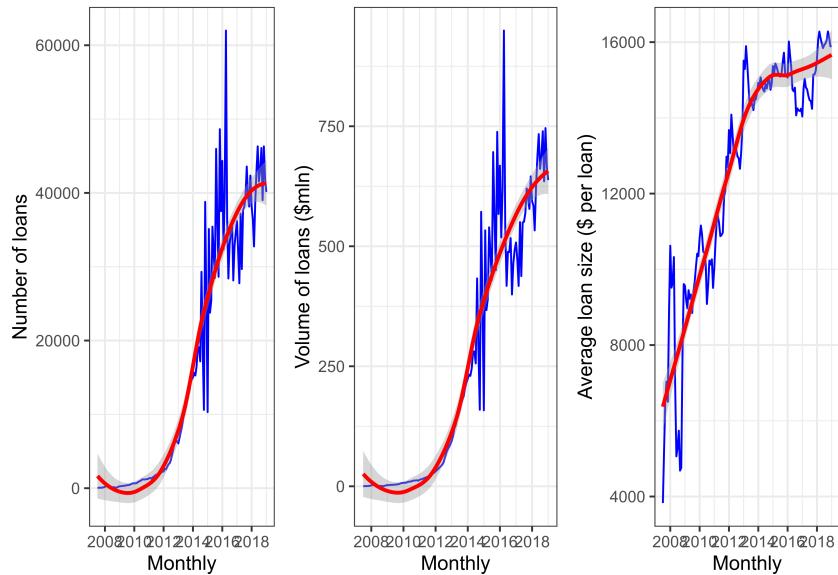


Figure 2.1: Business volume written per month

" 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 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 RHS Figure 2.1, on a sample of 100,000 loans). Partly, the increase in the early years could be interpreted success in improving marketing, distribution capabilities and confidence building. This metric plateau-ed

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

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, those metrics have essentially recovered to their 2016 level.

2.2.2 Loan lifecycle 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 possibility 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) from 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 Markov Chains transitions). 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.

²See LendingClub FAQ at [<https://help.lendingclub.com/hc/en-us/articles/215488038>] and help page [<https://help.lendingclub.com/hc/en-us/articles/216127897-What-happens-when-a-loan-is-charged-off->]

2.2.3 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 delinquencies). The dataset is very confusing in that regard: in the case of the **joint** applicant, it is clear that such information relates to before the loan is approved . In the case of the **principal borrower** however, the variable descriptions could be read as being pre-approval information, or post-approval 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 website³ 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](#)).

2.3 Rates

2.3.1 IRR and required credit margins

Figure 2.2 shows the evolution of credit margins over time grouped by ratings. The plots are made with a random sample of 300,000 loans.

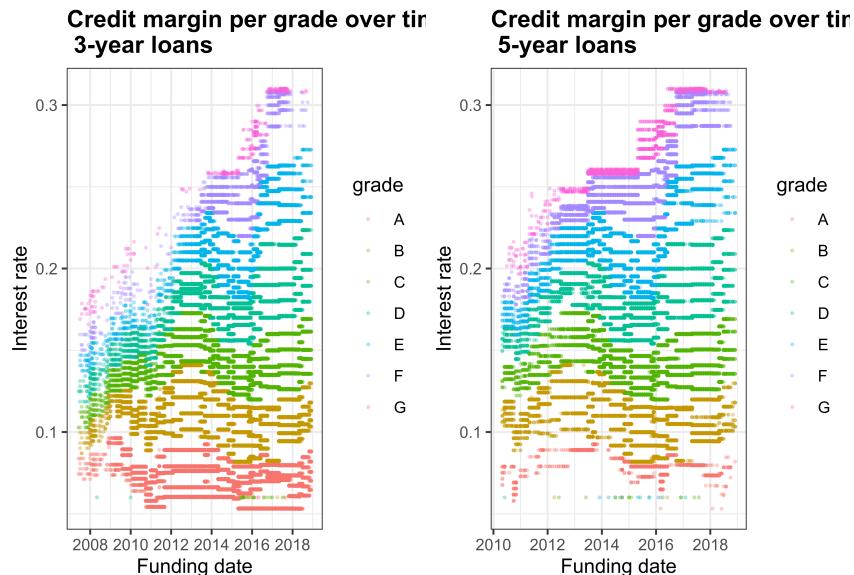


Figure 2.2: Credit margins per grade over time

³See <https://www.lendingclub.com/legal/prospectus>

We notice long periods where certain margins remain fairly stable which indicate that *both* the initial pricing was constant *and* that the proportion of default remains very low.

The graphs offer considerations that are relevant to the modeling:

- The margins clearly change over time. To the extent that they reflect a change in probability of default, the predictions will require to account for time (probably in a non-linear fashion).⁴
- For a given rating, the margins widen and narrow over time. The changes happen in multiples that depends on the ratings:
 - For high quality / low margin loans: the changes are multiples of the margin, for example going from roughly 3% to 6/7%.
 - Although the range of change is wide, those changes do not happen very often, especially in the later years.
 - By comparison, for low quality / high margin loans, the range of change is proportionally smaller, but more frequent and volatile.
- In other words, the relation between loan quality (its rating) and its pricing (the credit margin) will significantly non-linear.

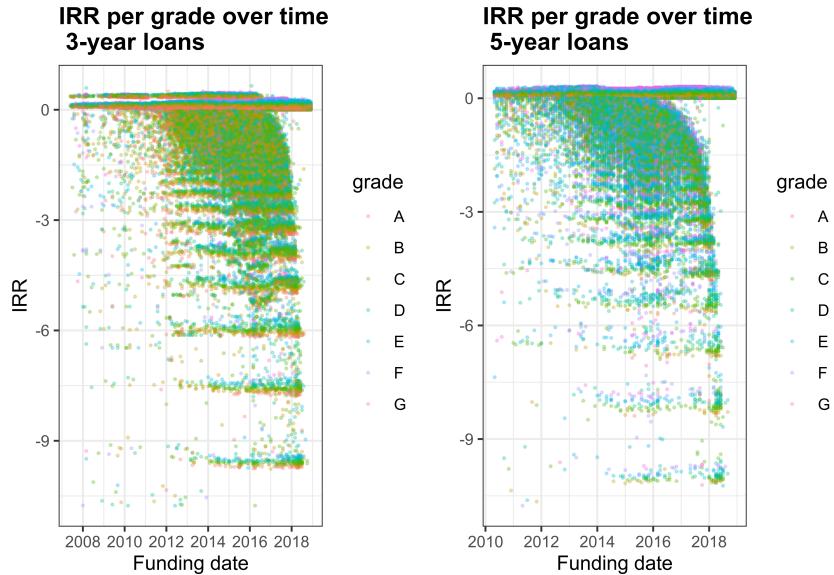


Figure 2.3: Credit margins per grade over time

2.3.2 Choice of predictors

Because we are interested decisions made prior to invest, we will limit the predictors to those that are realistically available prior to funding. We also remove information that is provided as a result of LC's own credit analysis (e.g. grade and interest rate).

2.3.3 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

⁴Note that we will add the second and third power of time (measure in months) to create this non-linearity. This will be the only real feature engineering that will be performed on the dataset.

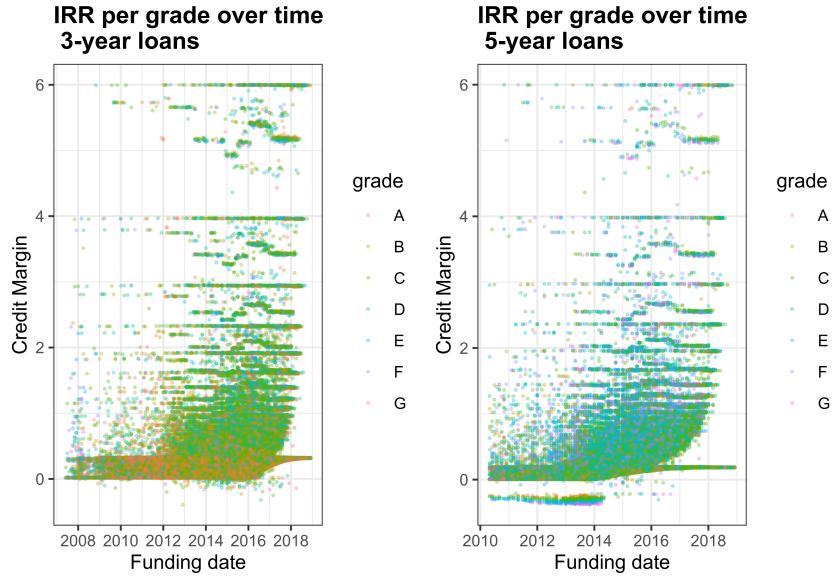


Figure 2.4: Credit margins per grade over time

determines the rate of interest according to a grid known to the borrower in advance⁵. The rates have changed over time. Those changes were not as frequent as market conditions (e.g. changes in Federal Reserve Bank's rates)⁶.

Figure 2.5⁷ shows the predetermined interest rate depending on the initial rating as of July 2019.

A1	6.46%	B1	10.33%	C1	14.30%	D1	18.62%
A2	7.02%	B2	11.02%	C2	15.24%	D2	20.55%
A3	7.56%	B3	11.71%	C3	16.12%	D3	23.05%
A4	8.19%	B4	12.40%	C4	16.95%	D4	25.65%
A5	8.81%	B5	13.08%	C5	17.74%	D5	28.80%

Figure 2.5: Interest rates given rating

At the date of this report, the ratings range from A (the best) down to D, each split in 5 sub-ratings. In the past, LC has also intermediated loans rated F or G (until 6 November 2017) and E (until 30 June 2019)⁸. 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 2.6 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 substantially widened over time. That is, the risk premium necessary to attract

⁵<https://www.lendingclub.com/investing/investor-education/interest-rates-and-fees>

⁶Corporate borrowers would negotiate interest margins on a case-by-case basis despite similar risk profiles.

⁷source: <https://www.lendingclub.com/investing/investor-education/interest-rates-and-fees>

⁸See <https://www.lendingclub.com/info/demand-and-credit-profile.action>

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 naturally considered safer (more A-rated, less G-rated). Identical ratings attract identical rates of interest.

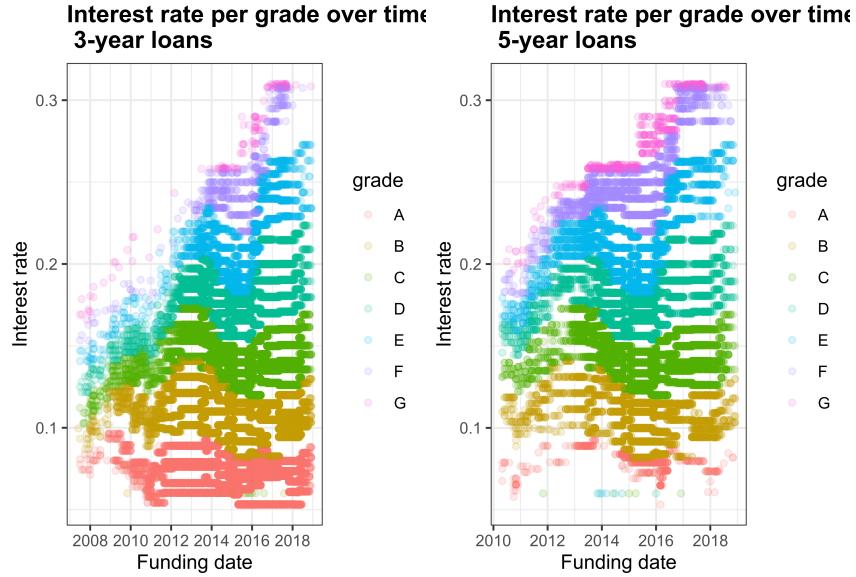


Figure 2.6: Interest rate per grade over time

By comparison, we plot the 3-year (in red) and 5-year (in blue) bank swap rates in Figure 2.7. 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 more relevant to examine the credit margins excluding swap rates.

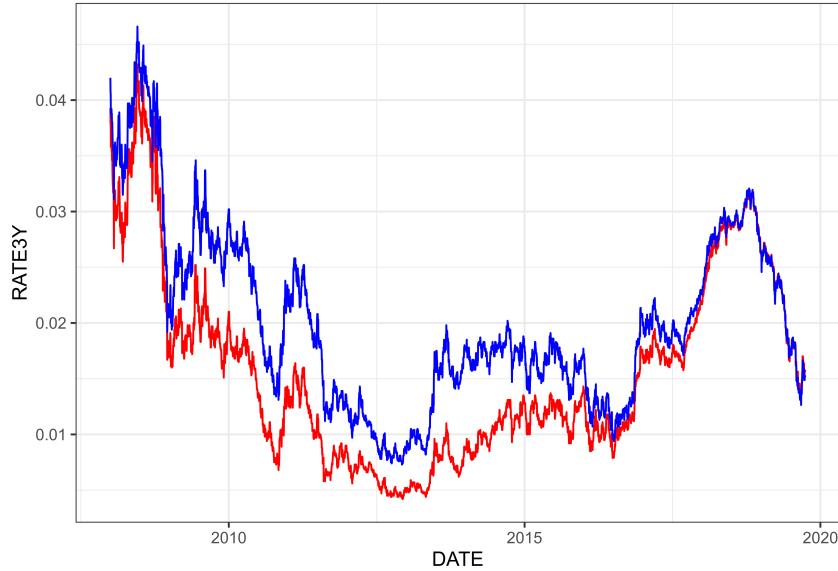


Figure 2.7: Historical Swap Rates

Figures 2.8 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.

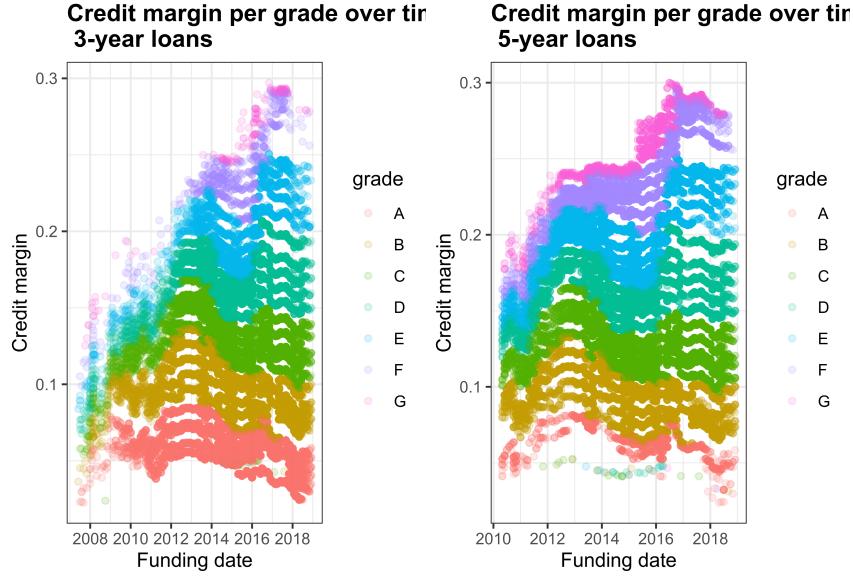


Figure 2.8: Credit margins per grade over time

2.3.4 Purpose

When applying, a potential borrower must state the purpose of the loan. As shown in table 2.2, by far the main purpose is the consolidation of existing debts.

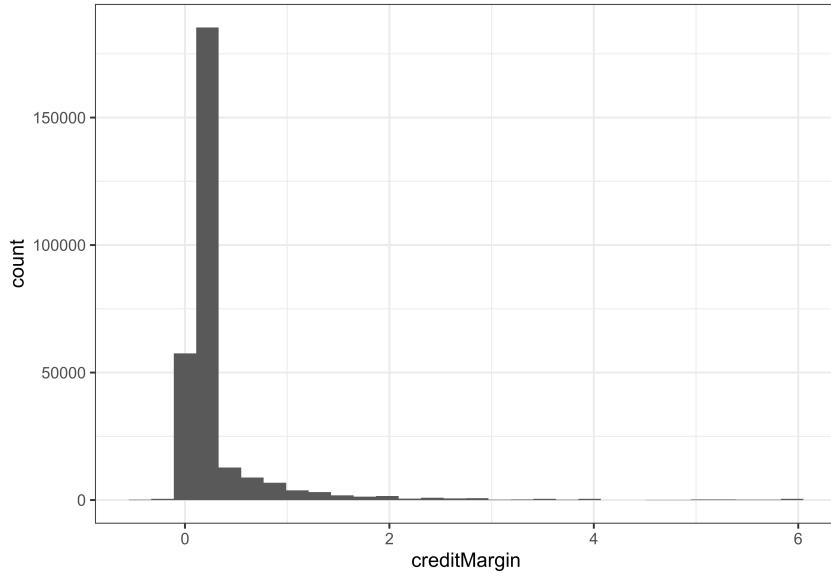


Figure 2.9: Histograms of credit margins per purpose

2.3.5 Payments

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

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

Where $Principal$ is the amount borrowed, $rate = \frac{\text{Quoted Interest Rate}}{12}$ is the monthly interest rate,

Table 2.2: Number of loans per purpose

Borrowing purpose	Count
debt_consolidation	758691
credit_card	286044
home_improvement	84709
other	75358
major_purchase	28451
small_business	15171
medical	15081
car	14184
moving	9218
vacation	8751
house	7011
wedding	2350
renewable_energy	914
educational	423

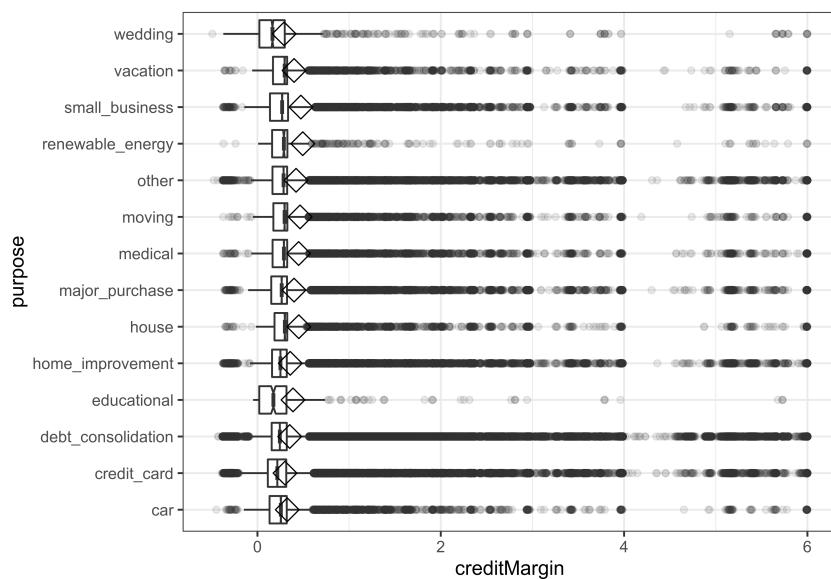


Figure 2.10: Boxplots of credit margins per purpose

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) ^ term), 2),
5       PMT_delta = abs(installment - PMT)
6     ) %>%
7     select(PMT_delta)
8
9
10  round(mean(100 * installmentError$PMT_delta), digits = 2)
11 })
```

2.4 Net present value

The behaviour of the NPV of loan losses is informative.

2.4.1 Average NPV and credit margin by subgrade

Figure 2.11 shows that as ratings worsen, the average NPV⁹ expressed as a portion of the funded amount decreases. For the best quality loans, we see that the NPV exceeds $1.00 = 100\%$: at a risk-free rate¹⁰, investors receive more than what is necessary to compensate for credit loss and can use the excess to cover additional costs mentioned in the Preamble. As ratings worsen, the NPV drops down to about 50%.

If loans were adequately priced, the excess returns (thanks to higher interest) should on average offset credit losses, that is an NPV average should be at least 100%. This seems to be the case down to ratings of about D4. Further down, credit losses become too frequent and/or too substantial to be covered on average. We posit that this justified rejecting loans applications rated E1 and below.

2.4.2 Distribution of principal losses by rating

We here consider nominal losses, that is not accounting for time effects (discounting).

Figure 2.12 shows that for a given grade, the losses are very widely spread. The loans are grouped by ratings and loans that have been fully repaid are removed.

Setting aside the loans rated “A” or “B”, the distributions seem log-normal. Unsurprisingly, the worse the rating the larger the principal loss.

2.4.3 NPV distribution by rating

Principal loss does not reflect the timing of that loss: for the same dollar amount, a loss now is worse than a loss later. This subsection looks at the NPVs of actual loan cashflow (principal and interest) discounted the risk-free rate.

Figure 2.13 shows that for a given grade, the NPVs are very widely spread. From top to bottom, loans are group by ratings: from quality ratings of A and B, average ratings of C and D, to

⁹The averages are *not* weighted by loan amount since an investor can invest in \$25 parcels. Weighting would have been appropriate if investors were instead forced to invest in the whole amount.

¹⁰Discounting a cash flow at the IRR gives a nil NPV. Discounting at a higher rate, resp. lower, gives a negative, resp. positive, NPV.

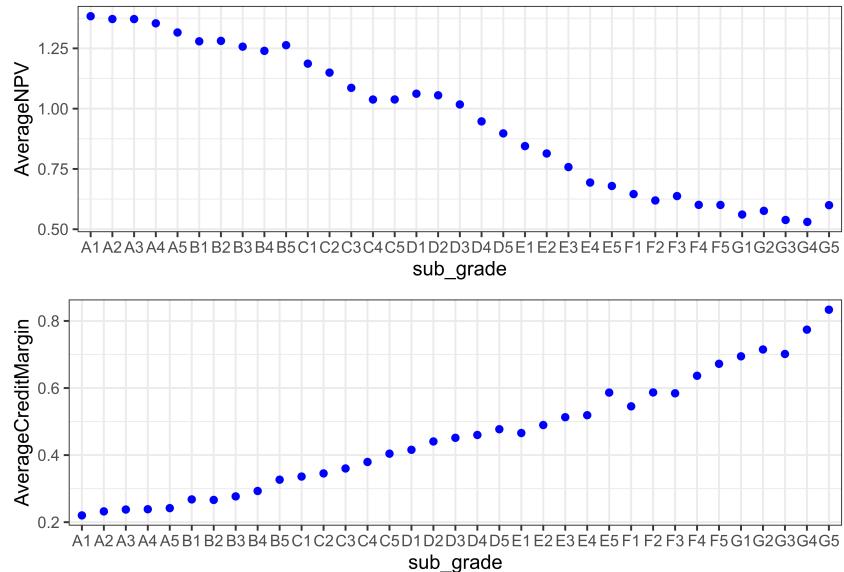


Figure 2.11: Average NPV et credit margin (%) depending on sub-rating

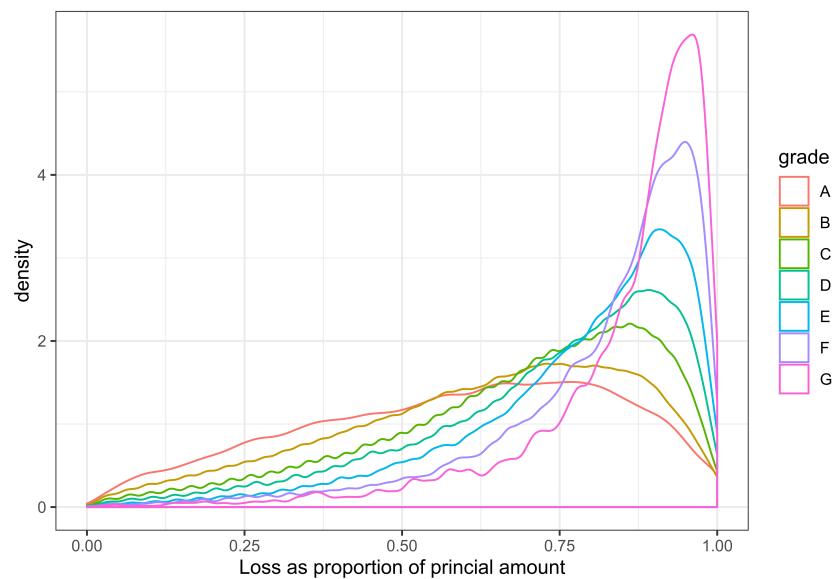


Figure 2.12: Distribution of the Principal Loss (%) depending on rating (y-axis square-root scaling)

poor ratings of E and below. From left to right, we focus on different parts of how NPVs are distributed. Note that each graph is based on a random sample of 100,000 loans (about 1/12th of the original set) and therefore the NPV densities are comparable from graph to graph. This spread is expected. The amount of loss for any loan can be anything a single dollar and 100%.

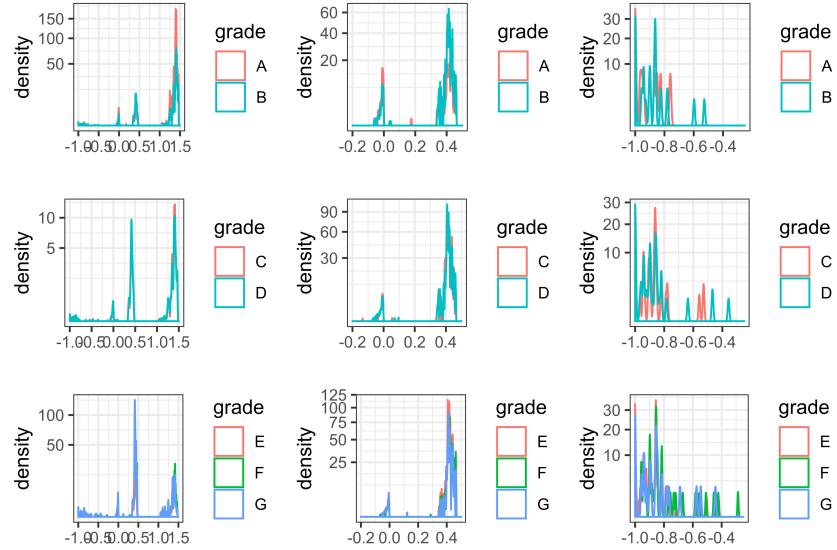


Figure 2.13: Distribution of NPV (%) depending on rating (y-axis square-root scaling)

At the outset, column by column (where NPVs are on the same scale), the NPV distribution show several modes on the same location. The modes are made more apparent by zooming on where the modes are present: the leftmost column basically shows the entire range of the NPVs (as portion of the loan). The middle graph zooms on the -20% / 50% range. The rightmost column zooms on the -100% / -25% section. Looking at the left hand scale, we can see that the lower NPVs overall gain in importance as the loan rating worsen.

Zooming without scaling the y-axis and grouping all the ratings available for investment on a single plot gives more details.

- Figure 2.14 shows a mode with a maximum around 1.25 / 1.5 being loans seemingly repaid in full (the mode is above 100% given the repayment of principal *and* interest);
- Figure 2.15 and figure 2.16 show a second and third mode around 41% and -1%;
- Finally, figure 2.17 one last very diffuse mode around -100%.

The overall trend is what we should expect. What is surprising is the existence of (1) very clearly defined modes which (2) are common to all types of borrowers. They roughly look log-normal, apart from the mode around 41% which look Gaussian.

2.5 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), and we will only retain loans which currently available (rated A1 to D5). From here on, everything will be based on this reduced dataset.

This reduced dataset contains 1306356 samples. Most of the loans (ca.80%) have been repaid in full (in other words **1 in 5 loans defaulted**). See Table 2.3.

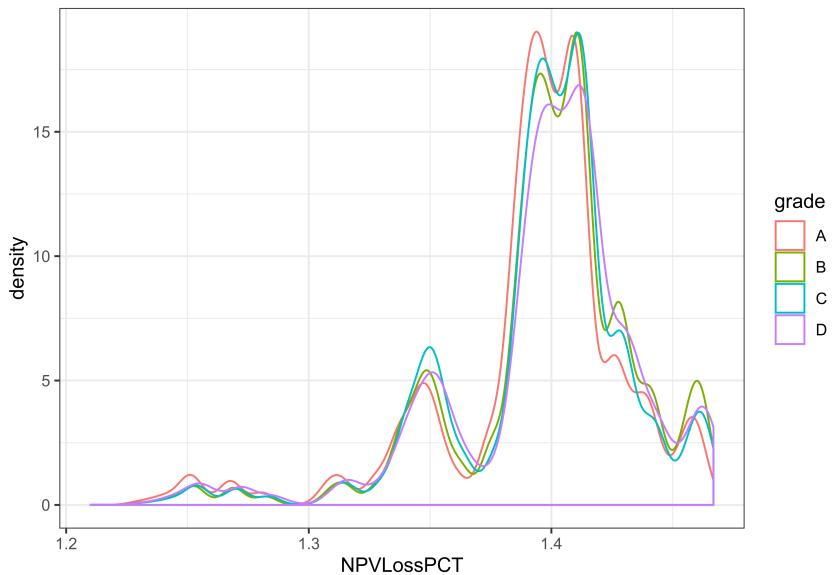


Figure 2.14: NPV % higher than 120% (no y-axis scaling)

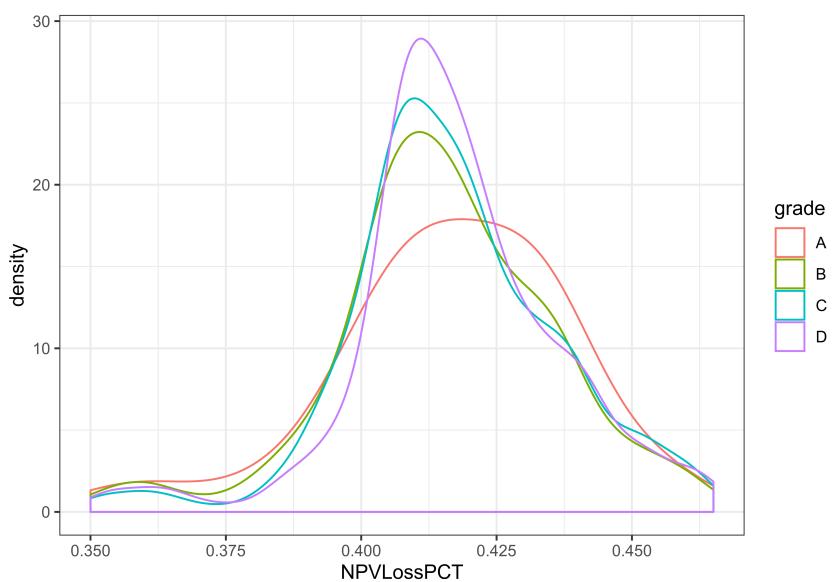


Figure 2.15: NPV % around 41% (no y-axis scaling)

Table 2.3: Matured loans per status

Loan status	Count	Proportion (%)
Fully Paid	964057	24101425
Charged Off	207546	5188650
Does not meet the credit policy.	1334	33350
Status:Fully Paid		
Does not meet the credit policy.	438	10950
Status:Charged Off		

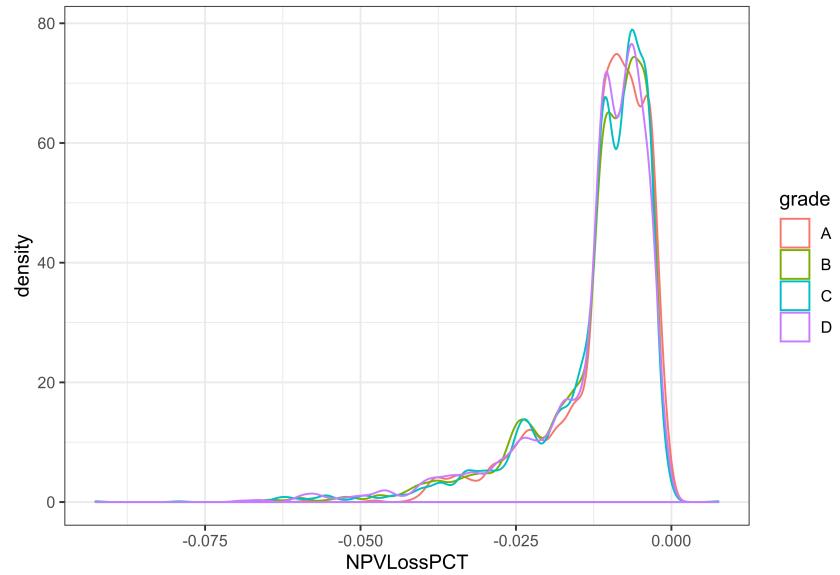


Figure 2.16: NPV % around -1% (no y-axis scaling)

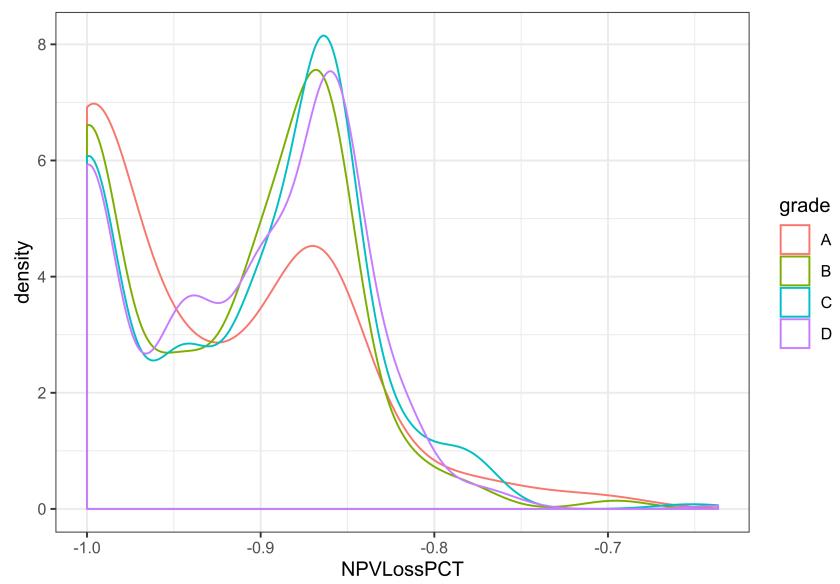


Figure 2.17: NPV % for close to total loss % (no y-axis scaling)

When grouped by grade (Figure 2.18), we see a clear correlation between grade and default: the lower the grade the higher the portion defaults (note the limited scale with a minimum at about 50%). In addition, in the reduced set most of the business is written in the B- or C-rating range.

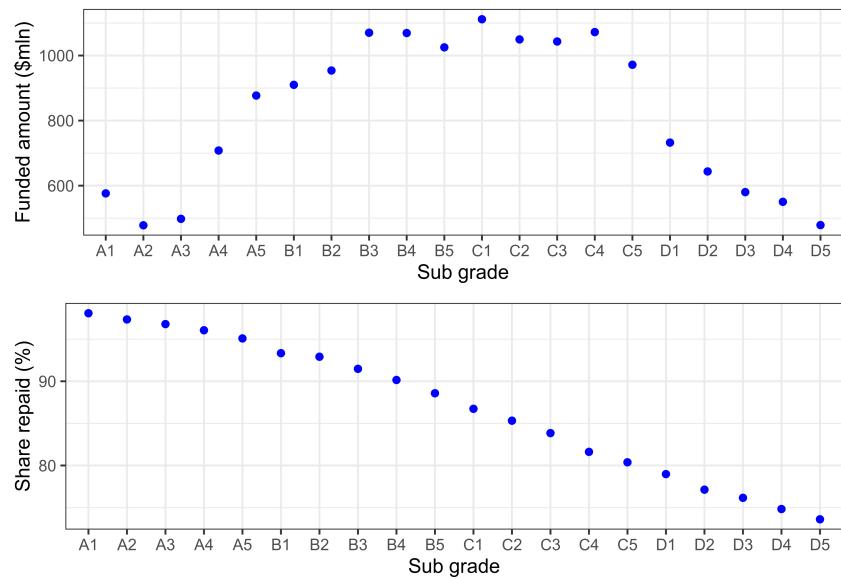


Figure 2.18: Funding and Write-offs by Sub-grades

Chapter 3

Logistic Regression Model and Credit Scorecard

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 number of predictors is potential large, in particular if we perform one-hot encoding of categorical values.
- The dataset has no context or reference point to interpret dollar amounts. We all intuitively understands that owing \$10 or \$1,000,000 are very different matters. Owing \$1,000 when living in New York is different from owing \$1,000 if living on \$2 per day in a developing country (surrounding economic environment matter). That intuition is shared economic knowledge, but that intuition is nowhere represented in the dataset. As readers, we automatically attribute that implicit knowledge to the data we read in the dataset. However, any model based on that data will never reflect that implicit knowledge if we do not supplement with external data. As shown in the previous section, credit margins have changed over time. This is clearly related to the wider US economic environment. Financial hardship is a key driver for some of the loans. Availability of disposable income is important to assess the ability to repay. Therefore, the cost of living, which varies from state to state, seems relevant.

WARNING:

IF YOU RUN THE MODELING SECTION, YOU WILL NEED UP TO 32GB OF MEMORY, AND EXPECT A LOT OF DISK SWAPPING WHEN DATASETS ARE JUGGLED IN AND OUT OF MEMORY TO DISK TO MINIMISE MEMORY USAGE. CALCULATION TIMES ARE BY BATCHES UP TO 10 MINUTES (DEPENDING ON HARDWARE). THIS EXCLUDES TIME NECESSARY TO PREPARE THE DATASET AS EXPLAINED ON THE PREVIOUS SECTIONS (THAT IS HOURS).

3.1 Logistic Regression

Logistic models (also called *logit model*) are used to model binary events. Examples would be passing or failing an exam, a newborn being a boy or a girl, a voter choosing one of two political parties, or – relevant to us – a borrower defaulting or not on a loan. If the binary variable is modeled as a 0/1 outcome, the model will yield a value between 0 and 1 which can be used as a probability.

We are interested in using a number of variables (being continuous and/or categorical) to model the binary response. A first port-of-call model is a linear combination of the variables. Since the predicted value would be continuous and not be bounded by 0/1, the outcome is transformed. A commonly used transformation of the logodds (logarithm of the odds given a particular probability) $\log\left(\frac{p}{1-p}\right)$. This expression has a few advantages: it converts any value (between $-\infty$ and $+\infty$ produced by the linear regression), and it is symmetrical around $x = 0$ and $y = 1/2$. That is, using the odds instead of the probability avoids infinity; it behaves identically when p approaches 0 or 1. The reciprocal of the logodds is $p = \frac{1}{1-e^{-x}}$.

For a number of X_i variables, the model to fit is then:

$$\text{logodds}(p) = \sum_i \alpha_i X_i \equiv p = \frac{1}{1 - e^{-\sum_i \alpha_i X_i}}$$

A commonly used format to evaluate the creditworthiness of a borrower is to create scorecards whereby particular characteristics are segmented into intervals and attributed discrete scores. In plain English, a continuous variable (say the age of the applicant) is segmented into intervals, called *bins*, (e.g. 0-18 year-olds, 18-26 year-olds, . . .), and then given a score. Those different segments become categorical variables. The task of the model is to:

- identify the best way to segment a continuous variable to maximise the information value of the different bins. Intuitively, empty or quasi-empty bins (either no or few applicants in the bin) are not very informative; bins for which the response is completely random are not informative, whereas a bin where the response always has the same response is informative (i.e. anybody with income of 0 and 10 dollars per year will default, anybody with a salary of \$1 million per month will repay).
- use a generalised linear model using the new categorical variables;
- transforms the linear coefficients estimated for each category into numerical scores.

3.1.1 Data preparation

The initial preparation of the dataset takes place in the `CleanLoan.Rmd` file where we bind:

- the raw LendingClub dataset (see file `Scripts/R-data-preparation.Rmd` for how it was prepared);
- from that raw set, only retain the variables that we selected to be used in the model (see column “Used in model?” in subsection 6.3);
- the NPV, IRR and credit margins calculated; and,
- the percentage of principal lost.

Some variables (such as relating to financial outcomes) are used for visualisations, not predictions.

```

1 # Quirk in bookdown?...
2 library(tidyverse)
3
4 #####
5 ##
6 ## select the variables that might be used to create the training+test set
7 ##
8
9 modelVarsIn <- c(LC_variable[LC_variable$inModel == TRUE, "variable_name"])$variable_name
10 modelVarsIn <- c(modelVarsIn,
11                     "grade_num", "sub_grade_num",
12                     "principal_loss_pct", "creditMargin", "monthDefault")
13
14 # Make sure that some variables are NOT in included in the final training set
15 modelVarsOut <- c("grade_num", "sub_grade_num",
16                   "principal_loss_pct", "creditMargin", "monthDefault",
17                   "zip_code")

```

We prepare a dataset with ONLY the predictors NOT removing NA's.

```

1 #####
2 ##
3 ## Prepare a dataset with ONLY the predictors NOT removing NA's
4 ##
5
6 loansPredictors <-
7   loansWorkingSet %>%
8
9   # Keep the chosen predictors
10  # Use tidyselect::one_of() to avoid errors if column does not exist
11  select(one_of(modelVarsIn)) %>%
12
13  ##
14  ## Dates to numeric, in 'decimal' years since 2000
15  ##
16  mutate_at(c("issue_d", "earliest_cr_line"), function(d) {
17    return(year(d) - 2000 + (month(d) - 1) / 12)
18  }) %>%
19
20  ## Add polynomials of the dates to model the time-trend shape
21  mutate(
22    issue_d2 = issue_d^2,
23    issue_d3 = issue_d^3,
24    earliest_cr_line2 = earliest_cr_line^2,
25    earliest_cr_line3 = earliest_cr_line^3
26  ) %>%
27
28  ## Create a logical flag TRUE for non-defaulted (good) loans
29  mutate(isGoodLoan = (principal_loss_pct < 0.001)) %>%
30
31  select(-tidyselect::one_of(modelVarsOut))

```

We split the dataset into a training (80%) and test set (20%).

```
1 ## #####  
2 ##  
3 ## Create training / test sets 80%/20%  
4 ##  
5 proportionTraining <- 0.8  
6 set.seed(42)  
7  
8 nSamples <- nrow(loansPredictors)  
9  
10 sampleTraining <- sample(1:nSamples, floor(nSamples * proportionTraining),  
11 replace = FALSE)  
12 loansTraining <- loansPredictors %>% slice(sampleTraining)  
13 loansTest <- loansPredictors %>% slice(-sampleTraining)  
14  
15 # Subsets of the training set  
16 set.seed(42)  
17 nSamplesTraining <- nrow(loansTraining)
```

We also create subsamples of the training set (20% thereof) for when quick calculation need making.

```
1 # 20%  
2 sample20 <- sample(1:nSamplesTraining, floor(nSamplesTraining * 0.20),  
3 replace = FALSE)  
4 loans20 <- loansTraining %>% slice(sample20)
```

3.2 Binning and Weight of Evidence

This subsection owes a debt to the source code of the `smbinning` package from which we reimplemented some aspects using the `tidyverse` style (the original source code uses SQL statements to access dataframes), and the documentation vignette of the `Information` package ¹. Our code is provided and imported as a package located on github ².

¹<https://cran.r-project.org/web/packages/Information/vignettes/Information-vignette.html>

²<https://github.com/Emmanuel-R8/SMBinning>

3.2.1 Background

Binning, Weight of Evidence (*WoE*) and Information Value (*IV*) have been widely used since the 1950's to convert continuous values to factors in a way that attempts to maximise the information content of the factors. This is achieved by adjusting the number and location of cut-off points to optimally partition the range of the continuous value.

After continuous variables are binned, all variables are categorical. Therefore WoE and IV measures then apply to either types of variables. They have some very important features:

WOE and IV enable to:

- Consider each variable's independent contribution to the outcome. This is $\mathcal{O}(n)$ for n variables instead of $\mathcal{O}(n^2)$ or $\mathcal{O}(n \log(n))$ of many algorithms;
- The WoE is a **factor-related** measure which assesses the relevance of each factor for a given variable;
- The IV is a **variable-related** measure which enables ranking variables between each other;
- The theoretical background to the measures lies in information theory. It cannot be over-emphasized that the measures **do not** use the values of the factors: it is a measure only calculated using the number of GOOD/BAD outcomes (see definitions below);
- Any NA values can be given their own factor: NAs are easily handled without filling values considerations. Given the previous point: the fact that NAs are present has no impact on the measures calculations! Sparse, incomplete or badly filled dataset are easily handled!
- The calculation is simple and quick (size of the dataset is not an issue in our case);
- The interpretation and visualisation of those measures is easy and intuitive;

3.2.2 Binning

The binning of a continuous variable is handled by the `partykit` package which produces *Conditional Inference Trees*. We will not describe the algorithm and the R package. See ([Hothorn et al., 2006](#)) for the theoretical background on Conditional Inference Trees, and ([Hothorn and Zeileis, 2015](#)) for a description of `partykit` implementation.

3.2.3 WOE and IV definitions

(This subsection is a modified copy for the `Information` R package vignette).

Let us have a binary dependent variable Y and a predictive variable X taking a discrete set of values x_1 to x_p (the factors). Y captures the loan defaults. Basically, the x_i represent one-hot encoding of the variable X .

In this situation, Naive Bayes can be formulated in a logarithmic form as:

$$\log \frac{P(Y = 1|x_1, \dots, x_p)}{P(Y = 0|x_1, \dots, x_p)} = \log \frac{P(Y = 1)}{P(Y = 0)} + \sum_{j=1}^p \log \frac{f(x_j|Y = 1)}{f(x_j|Y = 0)}$$

The naive Bayes model essentially says that the conditional log odds is equal to the sum of the individual factors (which will be the WoE). The word “naive” comes from the fact that this model relies on the assumption that all predictors are conditionally independent given Y , which is a highly optimistic (i.e. unrealistic) assumption.

3.2.3.1 Weight of Evidence

This can be remodeled in terms with $P(Y = 1|X = x_j)$ replaced by $GOOD_j$ being the proportion of good loans when only looking at bin j (how many good loans in that bin divided by the total number of loans in that bin), and similarly $P(Y = 0|X = x_j)$ replaced by BAD_j . We also define $GOOD$ as the proportion of good loans for the *entire variable* X . Here, proportion is interchangeable with probability.

We then have:

$$\log \frac{GOOD_j}{BAD_j} = \underbrace{\log \frac{GOOD}{BAD}}_{\text{sample log-odds}} + \underbrace{\log \frac{f(x_j|GOOD)}{f(x_j|BAD)}}_{\text{WoE}}$$

This relationship says that the conditional logit of $GOOD_j$ (odds of a good loan in a bin j), given x_j , can be written as the overall log-odds (total odds, i.e., the *intercept*) plus the log-density ratio – also known as the *Weight of Evidence*.

Note that the WoE and the conditional log odds of $Y = 1$ are perfectly correlated since the *intercept* is constant. Hence, the greater the value of WoE, the higher the chance of observing $Y = 1$. In fact, when WoE is positive the chance of observing $Y = 1$ is above average (for the sample), and vice versa when WoE is negative. When WoE equals 0 the odds are simply equal to the sample average.

3.2.3.2 Ties to Naive Bayes and Logistic Regression

Notice that the left-hand-side of the equation above – i.e., the conditional log odds of the variable – is exactly what we are trying to predict in a logistic regression model. Hence, when building a logistic regression model – which is perhaps the most widely used technique for building binary classifiers – we are actually trying to estimate the weight of evidence.

In our credit scoring situation, a “semi-naive” version of this model is quite popular. The idea is to transform the data into WoE vectors and then use logistic regression to fit the model

$$\log \frac{GOOD_j}{BAD_j} = \log \frac{GOOD}{BAD} + \sum_{j=1}^p \beta_j \log \frac{f(x_j|Y = 1)}{f(x_j|Y = 0)}$$

thus partly relaxing the assumption that all predictors in the model are independent (but not collinear). It should be noted that the underlying WoE vectors are still estimated univariately and that the coefficients merely function as scalars. For a more general model, GAM is a great choice.

As mentioned above, this relationship *does not* depend on the actual values of the bins; only the number of good and bad loans is used. If a bin represents NAs, the NAs have no impact on the calculation.

3.2.3.3 The Information Value

We can leverage WoE to measure the predictive strength of X – i.e., how well it helps us separate cases when $Y = 1$ from cases when $Y = 0$. This is done through the information value (IV) which is defined for continuous variables as:

$$IV_j = \int \log \frac{f(X_j|Y = 1)}{f(X_j|Y = 0)} (f(X_j|Y = 1) - f(X_j|Y = 0)) dx$$

Note that the IV is essentially a weighted “sum” of all the individual WoE values where the weights incorporate the absolute difference between the numerator and the denominator (WoE captures the relative difference).

3.2.3.4 Summary

More generally, now considering k variables, the most common approach to estimating the conditional densities needed to calculate WoE is to bin a variable X_k into individual bins $x_{k,j}$ and then use a histogram-type estimate.

$$\text{WoE}_{k,j} = \log \frac{\text{GOOD}_{k,j}}{\text{BAD}_{k,j}}$$

and the IV for variable X_k can be calculated as

$$\text{IV}_k = \sum_j (\text{GOOD}_{k,j} - \text{BAD}_{k,j}) \times \text{WoE}_{k,j}$$

3.2.4 Calculating WoE and IV

We took inspiration of the `smbinning` R package to optimally partition continuous variables into factors/bins. We however could not directly use this package as it does not interact with the tidyverse functions and uses SQL statements to access the content of dataframes. This solution enables to easily access SQL databases. But we note the the tidyverse takes the opposite approach of converting R statements into SQL statements.

We decided to implement a few functions that we will require (far from the entire `smbinning` package) as a new package called `binner` separately available on GitHub.

Let's install that package:

```

1 # Ensure that `binner` is here and available
2 if ("package:binner" %in% search()) {
3   detach("package:binner", unload = TRUE, force = TRUE)
4 }
5
6 if (!("binner" %in% installed.packages()[,1])) {
7   devtools::install_github("Emmanuel-R8/SMBinning")
8 }
```

and attach it to the environment.

```
1 library(binner)
```

3.2.5 Loop through all variables

Before creating the bins, we create a new dataframe and transform some string values to factors.

```
1 loansBinning <- loansTraining %>%
2   mutate(
3     home_ownership = as_factor(home_ownership),
4     emp_length = as_factor(emp_length),
5     grade = as_factor(grade)
6   )
```

For each variable we will create a new entry in a new tibble.

```
1 ######
2 ##
3 ## New tibble to store the list of bins + Weight of Evidences factors
4 ##
5 listBins <-
6   tibble(
7     variable = "",
8     type = "",
9     IV = 0.0,
10    WoE = list(),
11    .rows = 0
12  )
```

We then loop through each variable to create factors (for continuous variables) and calculate a table with the Weights of Evidence. We then calculate the Information Value.

```
1 # About 500 sec wall-time
2 startTime <- proc.time()
3
4 for (n in names(loansBinning)) {
5
6   # We don't test the response with itself
7   if (n != "isGoodLoan") {
8
9     # For categorical variable
10    if (class(loansBinning[[1, n]]) == "factor") {
11      # cat(" is a factor, ")
12      result <- WoETableCategorical(
13        df = loansBinning,
14        x = n,
15        y = "isGoodLoan",
```

```

16     maxCategories = 100)
17
18 } else {
19 # For continuous variable
20 result <- WoETableContinuous(df = loansBinning,
21                             x = n,
22                             y = "isGoodLoan",
23                             p = 0.05)
24 }
25
26 tryCatch({
27   if (is.na(result)) {
28     # In case no WoE table is create (of not enough bins)
29     add_row(
30       listBins,
31       variable = n,
32       type = NA,
33       IV = NA
34     )
35   } else {
36
37     listBins <- listBins %>%
38       add_row(
39         variable = n,
40         type = result$type,
41         IV = result$IV,
42         WoE = list(result$table)
43       )
44
45   }
46 },
47   finally = {})
48 }
49 }
50 cat("-- Lapsed time: ", proc.time() - startTime)
51 ## -- Lapsed time: 440.452 3.397 453.569 0 0

```

3.2.6 Select relevant variables

Table 3.1 guidelines are recommended to the relevance of variables given their Information Value (Bokhari, 2019).

The 15 best variables (in terms of Information Value in excess of 2%) are in Table 3.2:

Table 3.1: Variable relevance by Information Value

Information Value	Relevance
$IV < 2\%$	useless
$2\% < IV < 10\%$	weak
$10\% < IV < 30\%$	medium
$30\% < IV < 50\%$	strong
$50\% < IV$	warning: something is probably wrong!

Table 3.2: 15 top variables by IV

variable	IV
sub_grade	49.39
int_rate	46.53
grade	46.02
term	17.33
dti	7.45
verification_status	5.68
avg_cur_bal	5.43
tot_hi_cred_lim	5.00
num_tl_op_past_12m	4.85
mort_acc	4.38
total_bc_limit	4.12
issue_d2	4.09
issue_d3	4.09
issue_d	4.08
loan_amnt	3.62

```

1 bestBins <- listBins %>% filter(IV >= 0.02)
2
3 bestBins %>%
4   select(variable, IV) %>%
5   mutate(IV = round(100 * IV, digits = 2)) %>%
6   arrange(desc(IV)) %>%
7   slice(1:15) %>%
8   kable(caption = "15 top variables by IV", digits = 3)
```

However, we notice that this list includes variables that would not be available at the time the credit scoring is performed.

```

1 bestBins <-
2   bestBins %>%
3   filter(!(variable %in% c(
4     "loanID", "term", "int_rate", "creditMargin", "loan_status",
5     "grade", "sub_grade", "grade_num", "sub_grade_num",
6     "emp_length", "home_ownership", "monthDefault",
```

Table 3.3: 15 used top informative variables by IV

variable	IV
dti	0.074
verification_status	0.057
avg_cur_bal	0.054
tot_hi_cred_lim	0.050
num_tl_op_past_12m	0.049
mort_acc	0.044
total_bc_limit	0.041
issue_d2	0.041
issue_d3	0.041
issue_d	0.041
loan_amnt	0.036
mths_since_recent_inq	0.035
mo_sin_rcnt_tl	0.034
num_actv_rev_tl	0.034
open_rv_24m	0.033

```
7   "principal_loss_pct", "creditMargin", "monthDefault",
8   "isGoodLoan")))
```

The 15 most informative variables that will retain are in the following table. Interestingly, the square and cubic powers of the issue date are retained. (Recall subsection 2.3.1 on this being the only feature engineering.)

```
1 bestBins %>%
2   select(variable, IV) %>%
3   mutate(IV = round(IV, digits = 5)) %>%
4   arrange(desc(IV)) %>%
5   slice(1:15) %>%
6   kable(caption = "15 used top informative variables by IV", digits = 3)
```

And the 15 least informative (but retained) variables are:

```
1 bestBins %>%
2   select(variable, IV) %>%
3   mutate(IV = round(IV, digits = 5)) %>%
4   arrange(IV) %>%
5   slice(1:15) %>%
6   kable(caption = "15 least informative variables by IV", digits = 3)
```

Table 3.4: 15 least informative variables by IV

variable	IV
max_bal_bc	0.021
revol_util	0.025
open_rv_12m	0.026
inq_last_6mths	0.027
mo_sin_old_rev_tl_op	0.027
mths_since_recent_bc	0.027
mo_sin_rcnt_rev_tl_op	0.030
bc_open_to_buy	0.030
bc_util	0.030
num_rev_tl_bal_gt_0	0.032
percent_bc_gt_75	0.033
open_rv_24m	0.033
num_actv_rev_tl	0.034
mo_sin_rcnt_tl	0.034
mths_since_recent_inq	0.035

3.2.7 Create data table with one-hot encoding

Those variable will contain all the best characteristics. Every continuous variable is reformatted into factors reflecting the appropriate bins.

```

1 # Those variable will contain all the best characteristics. Every continuous variable is
2 # reformatted into factors reflecting the appropriate bins.
3 allFactorsAsCharacteristics <- loansTraining[, "loanID"]
4 allFactorsAsBins <- loansTraining[, "loanID"]
```

For each variable, we create new variables for each bin in the WoE table of that variable. Strictly speaking, this is not necessary for the generalised model algorithms. They are able to model using categories containing the factors. However, as we will see, the model will generate a number of NAs for variable (factors) which are co-linear. (This is the case for any linear model.) This will require removing those individual factors, which we found easier to do when each factor is given an individual variable (column).

```

1 for (index in 1:nrow(bestBins)) {
2   #for (index in 1:27) {
3     name <- bestBins$variable[index][[1]]
4
5     cat("---- Variable No. ",
6         index,
```

```

7     "## -- Name: ",
8     name, "\n")
9
10    ltIndex <- which(names(loansTraining) == name)
11
12    characteristic <- categoriseFromWoE(df = loansTraining[, ltIndex],
13                                         varName = name,
14                                         woeTable = bestBins$WoE[index][[1]])
15
16    bins <- categoriseFromWoE.Wide(df = loansTraining[, ltIndex],
17                                     varName = name,
18                                     woeTable = bestBins$WoE[index][[1]])
19
20    allFactorsAsCharacteristics <-
21    allFactorsAsCharacteristics %>%
22    cbind(characteristic)
23
24
25    allFactorsAsBins <-
26    allFactorsAsBins %>%
27    cbind(bins)
28
29 }

```

```

1 ## --- Variable No. 1 -- Name: loan_amnt
2 ## --- Variable No. 2 -- Name: verification_status
3 ## --- Variable No. 3 -- Name: issue_d
4 ## --- Variable No. 4 -- Name: dti
5 ## --- Variable No. 5 -- Name: inq_last_6mths
6 ## --- Variable No. 6 -- Name: revol_util
7 ## --- Variable No. 7 -- Name: open_rv_12m
8 ## --- Variable No. 8 -- Name: open_rv_24m
9 ## --- Variable No. 9 -- Name: max_bal_bc
10 ## --- Variable No. 10 -- Name: avg_cur_bal
11 ## --- Variable No. 11 -- Name: bc_open_to_buy
12 ## --- Variable No. 12 -- Name: bc_util
13 ## --- Variable No. 13 -- Name: mo_sin_old_rev_tl_op
14 ## --- Variable No. 14 -- Name: mo_sin_rcnt_rev_tl_op
15 ## --- Variable No. 15 -- Name: mo_sin_rcnt_tl
16 ## --- Variable No. 16 -- Name: mort_acc
17 ## --- Variable No. 17 -- Name: mths_since_recent_bc
18 ## --- Variable No. 18 -- Name: mths_since_recent_inq
19 ## --- Variable No. 19 -- Name: num_actv_rev_tl
20 ## --- Variable No. 20 -- Name: num_rev_tl_bal_gt_0
21 ## --- Variable No. 21 -- Name: num_tl_op_past_12m
22 ## --- Variable No. 22 -- Name: percent_bc_gt_75
23 ## --- Variable No. 23 -- Name: tot_hi_cred_lim
24 ## --- Variable No. 24 -- Name: total_bc_limit
25 ## --- Variable No. 25 -- Name: issue_d2
26 ## --- Variable No. 26 -- Name: issue_d3

```

Note that only 26 variables have been eventually retained out of the 145 of the initial dataset.

3.2.8 Comparison of individual characteristics

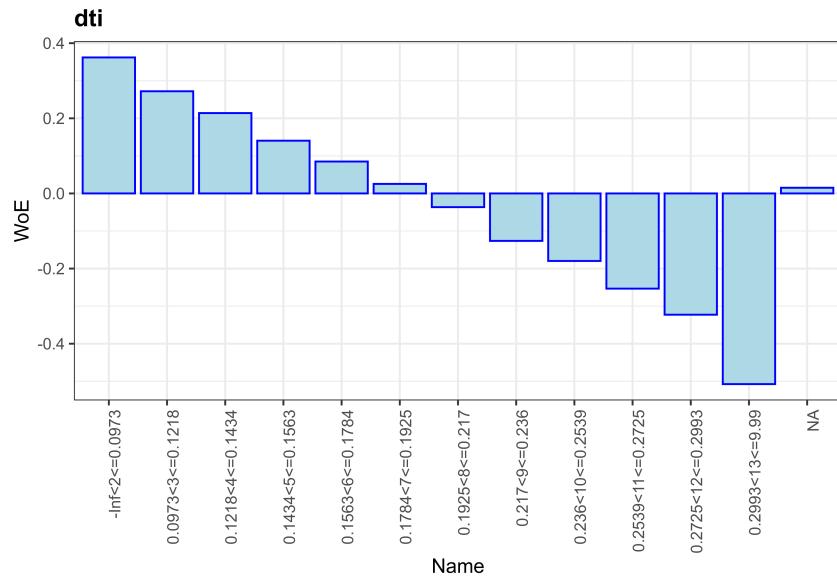
The following plots the 10 top variables weight of evidence plots. They show for those variables the weight of evidence of each individual factor. A positive WoE shows that the factor is positive to explain a positive outcome (that is that a loan is good).

The description of each variable is in the Appendix.

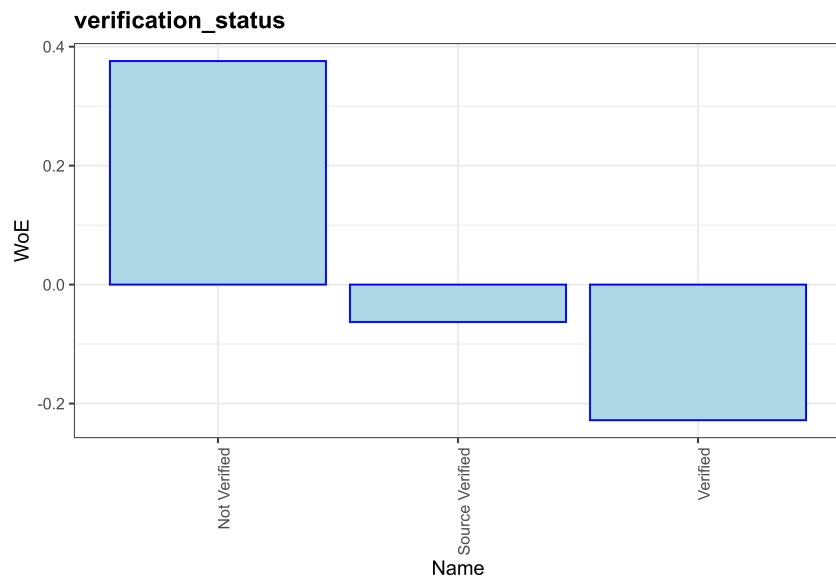
```

1 best10 <- bestBins %>% arrange(desc(IV)) %>% slice(1:10)
2
3 plotBinWoE <- function(n = 1) {
4   vName <- best10[n,]$variable
5
6   if (best10[n,]$type == "numeric") {
7     best10[n,]$WoE[[1]] %>%
8       arrange(WoE) %>%
9       ggplot(aes(Name, WoE)) +
10      geom_col(col = "blue", fill = "lightblue") +
11      theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
12      ggtitle(vName)
13   } else {
14     best10[n,]$WoE[[1]] %>%
15       arrange(WoE) %>%
16       ggplot(aes(Name, WoE)) +
17       geom_col(col = "blue", fill = "lightblue") +
18       theme(axis.text.x = element_text(angle = 90, hjust = 1)) +
19       ggtitle(vName)
20   }
21 }
22
23 }
24
25 listPlots <- lapply(1:10, function(n) { plotBinWoE(n) })
26
27 listPlots[[1]]

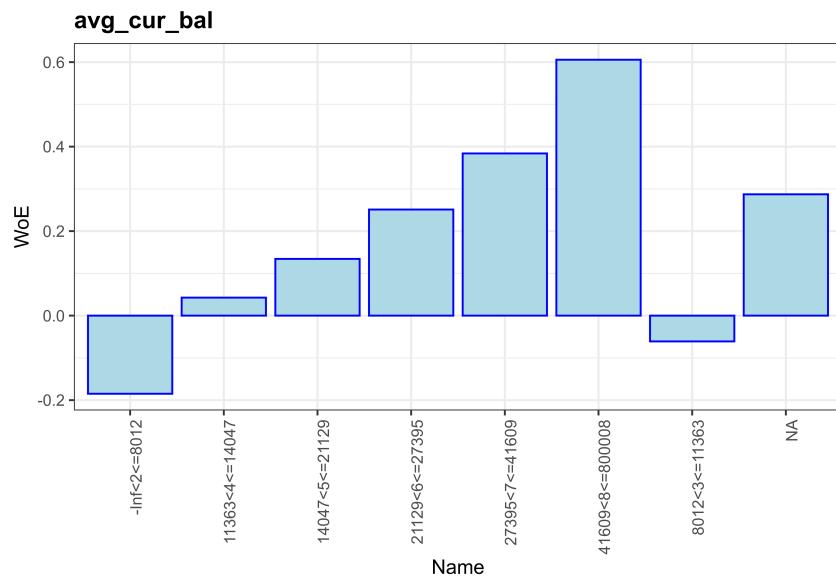
```



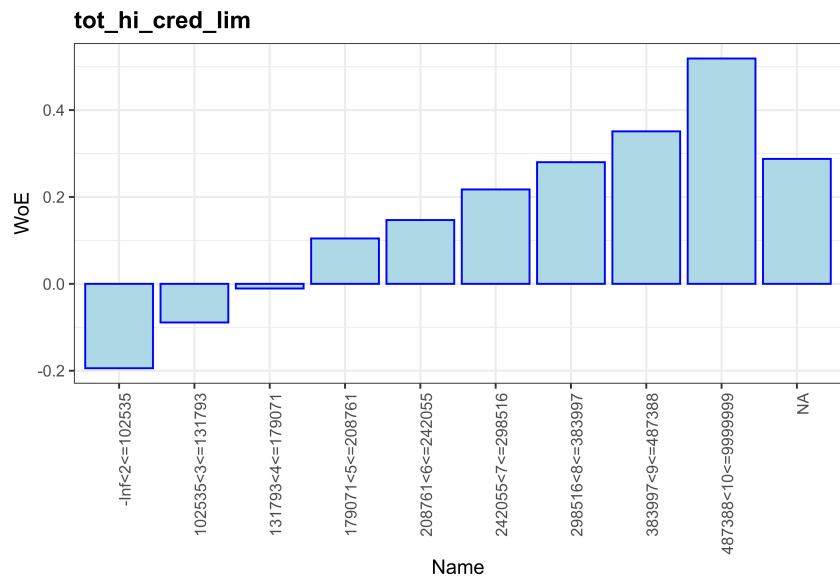
```
1 listPlots[[2]]
```



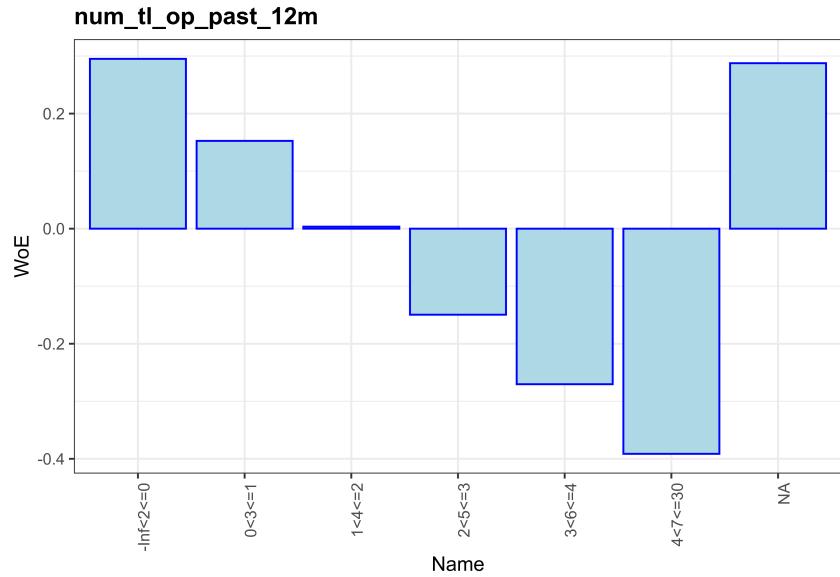
```
1 listPlots[[3]]
```



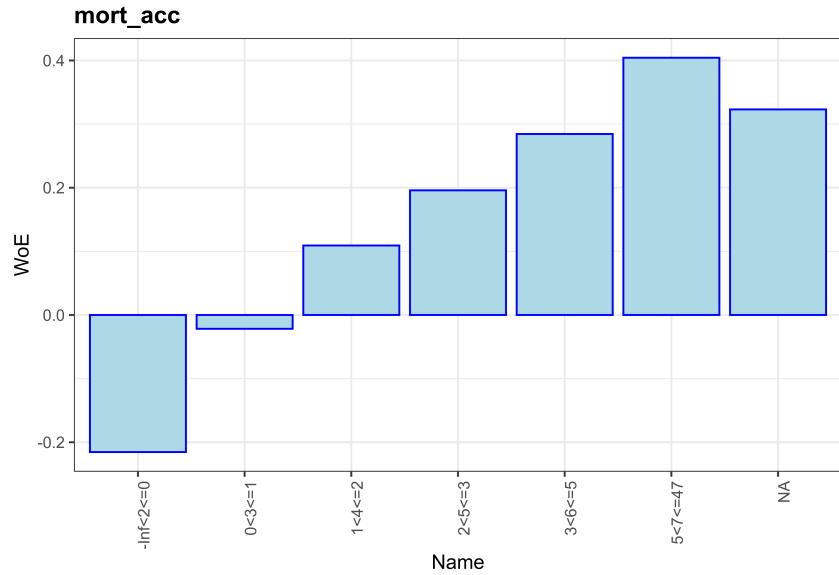
```
1 listPlots[[4]]
```



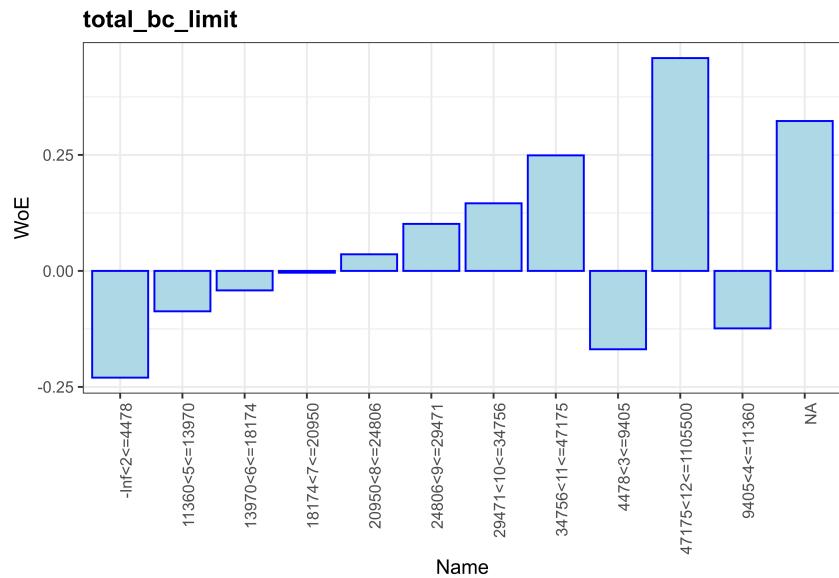
```
1 listPlots[[5]]
```



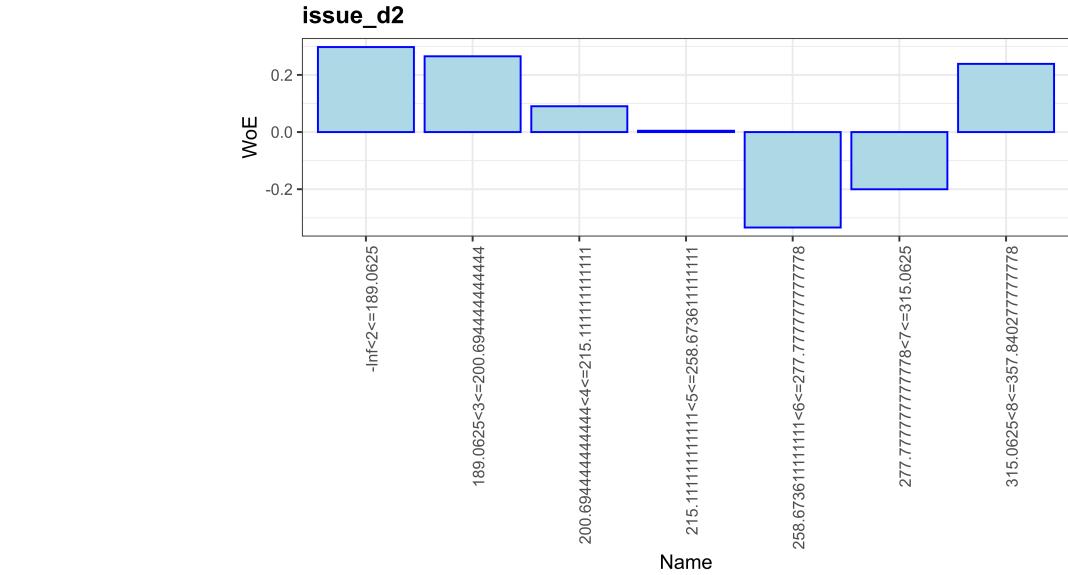
```
1 listPlots[[6]]
```



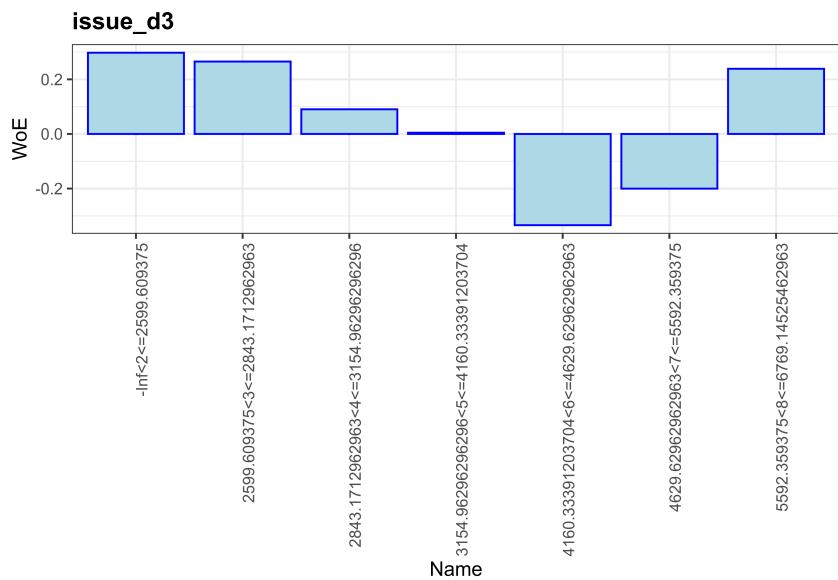
```
1 listPlots[[7]]
```



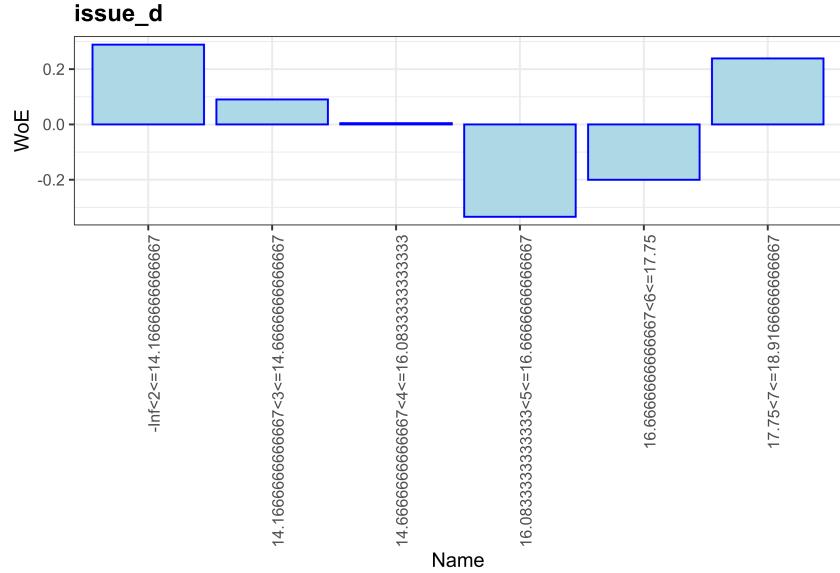
```
1 listPlots[[8]]
```



```
1 listPlots[[9]]
```



```
1 listPlots[[10]]
```



3.3 Logistic Regression

In this section, we fit a linear regression model using the fully binned dataset.

We tested several R linear regression packages. `glm` crashed on even small extracts of the dataset. `glmnet` returns errors that were not understandable nor documented on the internet. Although it does not include a `predict` function, we settled on the `speedglm` package which is designed to handle very large datasets. Note that it includes has a `select()` function which would shadow or conflict with `dplyr::select()`, so it is only used fully qualified to avoid any collision.

3.3.1 Remove identical bins

All linear regression algorithms require the variables to not have any linear relationship. Any collinearity prevents fitting a coefficient to such variables.

Although the original dataset may not contain such variables, binning variables into discrete factors will create such situations. The original variables were not identical, but when split into categories, some bins can become identical. For example, if a loan application does not provide income information (`income=NA` is TRUE for that loan), no debt-to-income ratio can be calculated (`dti=NA` will also be true for that loan). The `income` and `dti` variables were not collinear, but some of their bins can be.

To identify such bins, we use a “trick” to speed up comparison: we run a hash digest on each column and spot identical hashes.

```

1 namesCharacteristics <- names(loanSampleCharacteristics)
2 namesBins <- names(loanSampleBins)
3
4 # Starting list of names, calculate a hash value for each column with that name and
# identify duplicates.
5 duplicateNames <- names(
```

```

7   loanSampleBins[
8     duplicated(
9       lapply(loanSampleBins, digest::digest)))

```

We also identify any empty bin (this has only been useful when working on extremely small extracts of the original dataset).

```

1 # Remove any categories uniformly constant (only 0)
2 #
3 # This is important when working on a small extract of the dataset for variables that are
4 # seldom used (e.g. customers from certain states).
5 zeroColumnNames <- loanSampleBins[, colSums(loanSampleBins) == 0] %>% names()

```

We start with the binned dataset created in the previous subsection, removing variables with IV under 2%. All variables are one-hot encoded. Fitting the model will be done in four consecutive steps:

1. We fit a model on the binned dataset.
2. Next, we refit after removing any variable identified as colinear (i.e. NA coefficient).
3. Next, we refit after removing any variable whose significance is not at least p-value > 95%.

Finally, we will select a model.

3.3.2 First model - complete dataset

```

1 # About 700 sec wall-time to complete training dataset
2 # The dataset is the sample less duplicate, less zero-ed columns
3 {
4   doMC::registerDoMC(cores = 1) # Too many processes push over 32GB
5   startTime <- proc.time()
6
7   loansData <- loanSampleBins
8
9   SGLM_B_train <- speedglm::speedglm(isGoodLoan ~ .,
10                                         data = loansData,
11                                         family = binomial())
12
13  doMC::registerDoMC(cores = NULL)
14  cat("## Lapsed time: ", proc.time() - startTime)
15 }
1 ## -- Lapsed time: 169.877 10.025 200.529 0 0

```

There are a number of NA coefficients. The first 10 are:

```
1 speedglm:::summary.speedglm(SGLM_B_train)$coefficients %>%
2   as.data.frame() %>%
3   rownames_to_column(var = "Bin name") %>%
4   filter(is.na(Estimate)) %>%
5   slice(1:10) %>%
6   select("Bin name") %>%
7   kable(caption = "First model. Example of NA coefficients.", digits = 3) %>%
8   kableExtra::kable_styling(latex_options = "HOLD_position")
```

Table 3.5: First model. Example of NA coefficients.

Bin name
'(loan_amnt) 28000<loan_amnt<=40000'
'(verification_status) verification_status=Verified'
'(issue_d) 17.75<issue_d<=18.917'
'(dti) dti=NA'
'(inq_last_6mths) inq_last_6mths=NA'
'(revol_util) revol_util=NA'
'(open_rv_12m) open_rv_12m=NA'
'(open_rv_24m) 5<open_rv_24m<=53'
'(open_rv_24m) open_rv_24m=NA'
'(max_bal_bc) 7905<max_bal_bc<=776843'

The most significant bins are:

```
1 speedglm:::summary.speedglm(SGLM_B_train)$coefficients %>%
2   as.data.frame() %>%
3   rownames_to_column(var = "Bin_name") %>%
4   rename(z_value = "z_value") %>%
5   arrange(desc(z_value)) %>%
6   select(Bin_name, Estimate, z_value) %>%
7   slice(1:15) %>%
8   kable(caption = "First training. 15 best bins", digits = 3) %>%
9   kableExtra::kable_styling(latex_options = "HOLD_position")
```

Table 3.6: First training. 15 best bins

Bin_name	Estimate	z_value
'(loan_amnt) 3500<loan_amnt<=9000'	0.985	90.923
'(loan_amnt) -Inf<loan_amnt<=3500'	1.162	76.317
'(loan_amnt) 9000<loan_amnt<=10000'	0.738	59.375
'(loan_amnt) 10000<loan_amnt<=12000'	0.600	49.365
'(verification_status) verification_status=Not Verified'	0.284	40.127
'(loan_amnt) 14975<loan_amnt<=15000'	0.550	38.577
'(loan_amnt) 12000<loan_amnt<=14975'	0.504	38.228
'(num_tl_op_past_12m) 0<num_tl_op_past_12m<=1'	0.347	28.470
'(loan_amnt) 15000<loan_amnt<=17000'	0.376	27.665
'(mort_acc) 5<mort_acc<=47'	1.062	25.086
'(mort_acc) 3<mort_acc<=5'	1.013	24.616
'(mort_acc) 2<mort_acc<=3'	0.981	23.813
'(loan_amnt) 19950<loan_amnt<=28000'	0.239	23.742
'(mort_acc) 1<mort_acc<=2'	0.949	23.199
'(num_tl_op_past_12m) 1<num_tl_op_past_12m<=2'	0.252	22.480

We store the list of bins identified as colinear for later use.

```

1 NAsFirstTraining <-
2
3 # Take the results (just the estimated coefficients) from the model
4 tibble(
5   name = rownames(summary(SGLM_B_train)$coefficient),
6   estimate = summary(SGLM_B_train)$coefficient$Estimate
7 ) %>%
8
9 # Reformat the names to be the same as in the bins
10 mutate(name = stringr::str_remove_all(name, "\`")) %>%
11
12 # Make sure that response is not deleted by mistake and list all NAs
13 filter(name != "isGoodLoan" & is.na(estimate))
14
15 NAsFirstTraining <- NAsFirstTraining$name

```

3.3.3 Second training - fitted model less colinear bins

For the second training, we use the complete dataset with duplicate and colinear bins removed.

```

1 # Start the model training again
2 {
3   startTime <- proc.time()
4   doMC::registerDoMC(cores = 2) # Too many processes push over 32GB if more than 2
5
6   loansData <- loanSampleBins %>%
7     select(-one_of( c(duplicateNames,
8                   zeroColumnsNames,
9                   NAsFirstTraining)))
10
11  SGLM_B_retrain <- speedglm::speedglm(isGoodLoan ~ .,
12                                         data = loansData,
13                                         family = binomial())
14
15  doMC::registerDoMC(cores = NULL)
16  cat("-- Lapsed time: ", proc.time() - startTime)
17 }

1 ## -- Lapsed time: 153.658 16.549 234.563 0 0

```

The most significant bins are:

```

1 speedglm:::summary.speedglm(SGLM_B_retrain)$coefficients %>%
2   as.data.frame() %>%
3   rownames_to_column(var = "Bin_name") %>%
4   rename(z_value = "z value") %>%
5   arrange(desc(z_value)) %>%
6   select(Bin_name, Estimate, z_value) %>%
7   slice(1:15) %>%
8   kable(caption = "Second training. 15 best bins", digits = 3) %>%
9   kableExtra::kable_styling(latex_options = "HOLD_position")

```

Table 3.7: Second training. 15 best bins

Bin_name	Estimate	z_value
'(loan_amnt) 3500<loan_amnt<=9000'	0.985	90.923
'(loan_amnt) -Inf<loan_amnt<=3500'	1.162	76.317
'(loan_amnt) 9000<loan_amnt<=10000'	0.738	59.375
'(loan_amnt) 10000<loan_amnt<=12000'	0.600	49.365
'(verification_status) verification_status=Not Verified'	0.284	40.127
'(loan_amnt) 14975<loan_amnt<=15000'	0.550	38.577
'(loan_amnt) 12000<loan_amnt<=14975'	0.504	38.228
'(num_tl_op_past_12m) 0<num_tl_op_past_12m<=1'	0.347	28.470
'(loan_amnt) 15000<loan_amnt<=17000'	0.376	27.665
'(mort_acc) 5<mort_acc<=47'	1.062	25.086
'(mort_acc) 3<mort_acc<=5'	1.013	24.616
'(mort_acc) 2<mort_acc<=3'	0.981	23.813
'(loan_amnt) 19950<loan_amnt<=28000'	0.239	23.742
'(mort_acc) 1<mort_acc<=2'	0.949	23.199
'(num_tl_op_past_12m) 1<num_tl_op_past_12m<=2'	0.252	22.480

We can now select any bins whose significance value is under 2σ .

```
1 NAsSecondTraining <-
2   tibble(
3     name = rownames(summary(SGLM_B_retrain)$coefficient),
4     estimate = summary(SGLM_B_retrain)$coefficient$Estimate,
5     zValue = abs(summary(SGLM_B_retrain)$coefficient$"z value")
6   ) %>%
7
8   mutate(name = stringr::str_remove_all(name, "\`")) %>%
9   filter(name != "isGoodLoan") %>%
10
11  # Remove stray NAs or anything less than 2 sigmas
12  filter(is.na(estimate) | zValue < 2)
13
14 NAsSecondTraining <- NAsSecondTraining$name
```

3.3.4 Third training - second fitted model less non-significant bins

For the third training, we use the complete dataset with duplicate, colinear and non-significant bins removed.

```
1 # Start the model training again. About 350 sec.
2 {
3
4   startTime <- proc.time()
5   doMC::registerDoMC(cores = 2)
6
7   loansData <- loanSampleBins %>%
8     select(-one_of( c(
9       duplicateNames,
10      zeroColumnsNames,
11      NAsFirstTraining,
12      NAsSecondTraining
13    )))
14
15   SGLM_B_reretrain <- speedglm::speedglm(isGoodLoan ~ .,
16                                             data = loansData,
17                                             family = binomial())
18
19   doMC::registerDoMC(cores = NULL)
20   cat("-- Lapsed time: ", proc.time() - startTime)
21 }
22
23 ## -- Lapsed time: 54.173 5.597 77.389 0 0
```

The most significant bins are:

```

1 speedglm:::summary.speedglm(SGLM_B_reretrain)$coefficients %>%
2   as.data.frame() %>%
3   rownames_to_column(var = "Bin_name") %>%
4   rename(z_value = "z value") %>%
5   arrange(desc(z_value)) %>%
6   select(Bin_name, Estimate, z_value) %>%
7   slice(1:15) %>%
8   kable(caption = "Third training: 15 best bins", digits = 3) %>%
9   kableExtra::kable_styling(latex_options = "HOLD_position")

```

Table 3.8: Third training: 15 best bins

Bin_name	Estimate	z_value
'(loan_amnt) 3500<loan_amnt<=9000'	0.938	87.483
'(loan_amnt) -Inf<loan_amnt<=3500'	1.108	73.272
'(loan_amnt) 9000<loan_amnt<=10000'	0.703	56.953
'(loan_amnt) 10000<loan_amnt<=12000'	0.562	46.561
'(verification_status) verification_status=Not Verified'	0.323	45.961
'(num_tl_op_past_12m) -Inf<num_tl_op_past_12m<=0'	0.570	42.824
'(num_tl_op_past_12m) 0<num_tl_op_past_12m<=1'	0.462	41.164
'(loan_amnt) 14975<loan_amnt<=15000'	0.530	37.341
'(num_rev_tl_bal_gt_0) -Inf<num_rev_tl_bal_gt_0<=2'	0.481	37.093
'(loan_amnt) 12000<loan_amnt<=14975'	0.458	35.041
(Intercept)	0.640	32.243
'(mort_acc) 5<mort_acc<=47'	1.169	31.135
'(num_tl_op_past_12m) 1<num_tl_op_past_12m<=2'	0.329	31.083
'(mort_acc) 3<mort_acc<=5'	1.106	30.336
'(num_rev_tl_bal_gt_0) 4<num_rev_tl_bal_gt_0<=5'	0.328	29.163

Table 3.9 shows the Akaike Information Criterion and Log-likelihood for each of the 3 training results. Two points to note:

- The results did not change from first to second training. This is normal since we only removed non-informative categorical bins for which the regression could not determine a coefficient (i.e. collinearity with other variables). We did not remove any information;

- Removing seemingly non-significant bins makes the model worse (at least by those measures). This would suggest that bins should not be individually considered, but only the entire variable to which they belong.

Table 3.9: Training AIC and Log-likelihood

Criteria	Training #1	Training #2	Training #3
AIC	977037.2	977037.2	982957.0
Log likelihood	-488350.6	-488350.6	-491370.5

In the next section we will only work with the second model since it is the best (by those measures).

3.4 Model results

3.4.1 Final list of variables

We selected the model fitted during the second training.

```
1 # Model to use
2 GLModel <- SGLM_B_retrain
```

We collect the list of selected bins.

```
1 # List of model variables
2 modelNames <-
3   attr(GLModel$coefficients, "names") %>%
4   enframe(x = .) %>%
5   rename(variableName = "value") %>%
6   select(variableName) %>%
7   mutate(variableName = str_remove_all(variableName, "\`"))
```

And we collect the model summary (coefficients, significance, etc.).

```

1 # and their coefficients in the model (the summary function for speedglm objects is not exported and
2 # bookdown seems to have a problem with that).
3
4 GLMCoefficients <-
5   speedglm:::summary.speedglm(GLModel)$coefficients %>%
6     as_tibble() %>%
7     cbind(modelNames) %>%
8     rename(
9       zValue = "z value",
10      pValue = "Pr(>|z|)",
11      stdErr = "Std. Error"
12    ) %>%
13
14 # reorder columns to have names first
15 select(variableName, everything())

```

3.4.2 Scoring

Scoring expresses the coefficients that were estimated during the logistic regression into points on a scale. The model estimates the log-odds that a loan will not default. We will provide a detailed description of the calculations to, hopefully, provide guidance to others. (The references we have found in the course of our research have been surprisingly limited.)

Let us recall that our model will be in the form of a linear regression to model the logodds of the probability of default. That is:

$$\text{logodds}(p) = \text{intercept} + \sum_{\text{Variable}=k} \sum_{\text{bin}=i} \alpha_{k,i} x_{k,i}$$

where:

$$\text{logodds}(p) = \log \frac{p}{1-p}$$

Thanks to the properties of the logodds transformation, the model can take any value between $-\infty$ and $+\infty$ and yield valid probability values.

The idea of scoring is to replace this linear relationship with a points system: if an applicant ticks a box on a particular question/variable, he/she gets so many points. For example, asked for an age band, an applicant would receive 10 points if between 18 and 25 year-old, and 25 points if between 25 and 32. The age bands (the variable bins) were calculated above.

Basic convenience and marketing common sense dictates that:

- points should be rounded (who wants to see 14.99432529875 on an application form?)
- those points should not be negative.
- Scorecards do not start with an initial piggybank of points. In terms of modeling, that means no intercept.

The idea of scorecard applies to model similar to this one. Marketing consideration might be relevant to our context. They would be completely irrelevant to scorecard for a disease risk assessment or probability of failure of a mechanical part.

Scoring will replace the $\alpha_{k,i}$ coefficients with scores $S_{k,i}$. In addition, we will get rid of the intercept. To remove the intercept, we will apportion it equally across all variables:

$$\text{logodds}(p) = \sum_{\text{Variable}=k} \left(\frac{\text{intercept}}{K} + \sum_{\text{bin}=i} \alpha_{k,i} x_{k,i} \right)$$

where K is the number of variables. And if I_k is the number of bins for variable k :

$$\text{logodds}(p) = \sum_{\text{Variable } k=1}^K \sum_{\text{bin } i=1}^{I_k} \left(\frac{\text{intercept}}{K \cdot I_k} + \alpha_{k,i} x_{k,i} \right)$$

Then we select a factor F that will dilate/contract the coefficients such that:

$$\text{logodds}(p) = \sum_{\text{Variable}=k} \left(\frac{\text{intercept}}{K} + \sum_{\text{bin}=i} \alpha_{k,i} x_{k,i} \right)$$

3.4.3 Model scorecard

The conversion is done using three parameters that are chosen somewhat arbitrarily and which define a line:

- the number of points increase / decrease that would reflect halving / doubling the odds of defaulting;
- an *anchoring* score reflecting a particular odd.

For our purpose, we will choose 5,000 points ($\text{Score}_{\text{anchor}}$) being equivalent to 1 in a 20 to default ($\text{Odds}_{\text{anchor}}$), i.e. $\text{Score}_{\text{anchor}} 5,000 \text{ points} \Leftrightarrow \text{Odds}_{\text{Anchor}} = \frac{1/20}{1-1/20}$. We will also choose 100 to reflect *times 2* change in odds (*DoubleOdds*). Those choices are completely arbitrary. Basically, the score is a linear representation of the odds. The score is defined by a point (the anchor) and the slope of the line going through that point.

Those values will reflect the total estimated score for a borrower (i.e. loan sample). The number of characteristics (information points gathered in a credit application), or number of bins, have to be irrelevant in calculating this score. In other words, if LendingClub were to gather 5 more information points, the score should, **mutatis mutandis**, be unchanged. (However, we would hope that the quality of the estimated score would improve.)

The score per variable needs to be adjusted using the number of information points, that is the number of characteristics. Here, the model has been trained on the number of bins. We first need to determine how many characteristics are used in the model.

```

1 # The list of characteristics is extracted from the list of bins names.
2 numberOfRowsCharacteristics <-
3
4 modelNames %>%
5
6 # List of all names excluding the intercept which is not part of the scoring calculations (it
7 # would mean giving points for free as a base line)
8 filter(variableName != "(Intercept)") %>%
9
10 # Extract strings ("[a-zA-Z0-9_]*) preceded by an opening bracket ("(?<=\\"()

```

```

11  # closing bracket ("(?=\))"). (The column names were formatted this way for that purpose.) See
12  # "https://github.com/rstudio/cheatsheets/blob/master/strings.pdf" for details on the regex.
13  mutate(characteristic = str_match(variableName, "(?<=\(\)[a-zA-Z0-9-_]*(?=\\\))") %>%
14
15  # We are only interested in how many distinct characteristic there are.
16  distinct(characteristic) %>%
17  nrow()
18
19  numberOfCharacteristics
20
21 ## [1] 24

```

We can now perform the scoring calculation.

```

1 ProbDefaultAtAnchor <- 1/20
2
3 # The model is trained so that the 'Good' outcome is that a loan does not default. The odds that
4 # therefore calculated on the probability of no default.
5 ProbAtAnchor <- 1 - ProbDefaultAtAnchor
6 OddsAnchor <- ProbAtAnchor / (1 - ProbAtAnchor)
7 ScoreAnchor <- 2000
8
9 # Doubling odds = 100 points
10 DoubleOdds <- 100
11
12 ScoreFactor <- OddsAnchor / log(2)
13
14 # 5,000 points is 20:1 odds of default
15 ScoreOffset <- ScoreAnchor - ScoreFactor * log(OddsAnchor)

```

Then we apportion across characteristics.

```

1 # Score at the intercept
2 Intercept <- summary(GLModel)$coefficients["(Intercept)", "Estimate"]
3
4 ScorePerVariable <- (ScoreFactor * Intercept + ScoreOffset) / numberOfCharacteristics
5 InterceptPerVariable <- Intercept * ScoreFactor / numberOfCharacteristics
6
7 GLMScores <-
8   GLMCoefficients %>%
9   mutate(
10     weight = Estimate * ScoreFactor,
11     weightScaled = weight + ScorePerVariable,
12     points = round(weightScaled)
13   )

```

Very important: The intercept points have been allocated across all characteristics. Therefore the regression coefficient estimated for the intercept becomes redundant and needs removing.

```
1 GLMScores[1, "points"] <- 0
```

The `speedglm` package does not have a `predict` function once models have been trained. However, using the model is a simple matrix multiplication: Loan matrix \times Scorecard weights

3.5 Training set

We first review how the fitted model performed on the training set.

We remove every variable that is not in the list of variables in the model then convert into a matrix.

```
1 # Remove every variable that is not in the list of variables in the model then convert
2 # into a matrix
3 allMatrix <-
4   allFactorsAsBins[, !is.na(match(
5     names(allFactorsAsBins),
6     str_remove_all(GLMCoefficients$variableName, "\`"))
7   )] %>%
8   as.matrix()
```

The coefficients of the model also include the intercept (set at zero in the scores). We add a column of 1's to the training data.

```
1 # Add a column of 1s for the intercept
2 allMatrix <-
3   cbind(as.vector(rep.int(
4     x = 1, times = dim(allMatrix)[1]
5   )), allMatrix)
6   dim(allMatrix)
1 ## [1] 1045084    168
```

The coefficients of the model, then the scores, are converted to a vector format.

```
1 CoefficientsVector <- GLMCoefficients$Estimate %>% as.matrix()
2
3 # Score per variable
4 TrainingScorecard <- allMatrix %*% ( GLMScores$points %>% as.matrix() )
```

We can now multiply the matrix of sample with the vector of coefficients.

```
1 TrainingLogit <- allMatrix %*% CoefficientsVector
2 TrainingLogit <-
3   enframe(TrainingLogit[, 1]) %>%
4   mutate(oddsGood = exp(value),
5         p = 1 / (1 + oddsGood)) %>%
6   cbind(TrainingScorecard)
```

3.5.1 Densities of the training results

We plot the results of the training model and group the results by rating (“A” to “G”) in Figure 3.1.

```
1 gridExtra::grid.arrange(
2   loansTraining %>%
3     cbind(TrainingLogit) %>%
4     filter(between(value, -2, 5)) %>%
5     ggplot(aes(value, col = grade)) +
6     geom_density(adjust = 0.5) +
7     ggtitle("Logit value"),
8
9   loansTraining %>%
10    cbind(TrainingLogit) %>%
11    ggplot(aes(p, col = grade)) +
12    geom_density(adjust = 0.5) +
13    scale_x_log10() +
14    ggtitle("Probability of being GOOD"),
15
16   loansTraining %>%
17    cbind(TrainingLogit) %>%
18    filter(between(oddsGood, 0.1, 100)) %>%
19    ggplot(aes(oddsGood, col = grade)) +
20    geom_density(adjust = 0.5) +
21    scale_x_log10() +
22    ggtitle("Odds of being GOOD"),
```

```

23
24   loansTraining %>%
25     cbind(TrainingLogit) %>%
26     ggplot(aes(TrainingScorecard, col = grade)) +
27     geom_density(adjust = 1) +
28     ggtitle("Scorecards"),
29
30   ncol = 2, nrow = 2
31 )

```

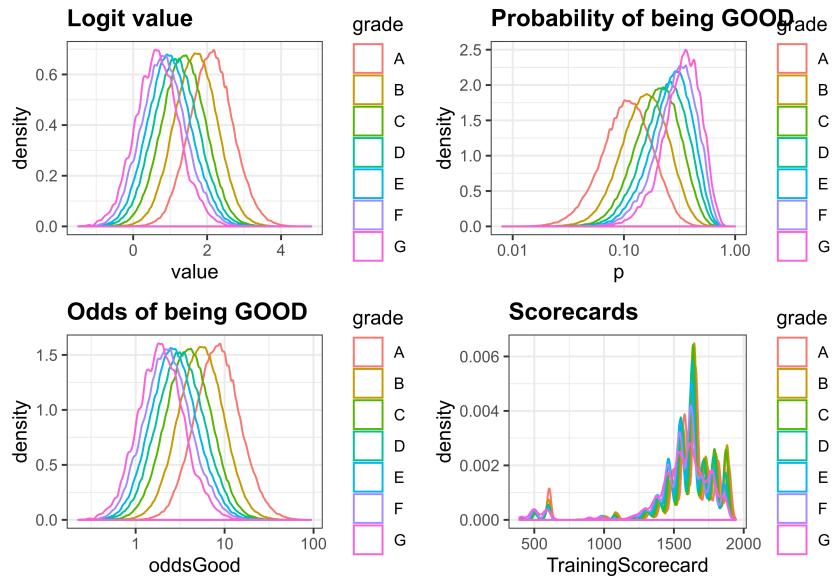


Figure 3.1: Model results on the training set

3.6 Test set

We replicate the exact same steps on the test set loans, also converted into a matrix.

```

1 predictionCategories <- loansTest[, "loanID"]
2
3 for (index in 1:length(bestBins$variable)) {
4   binned <-
5     binner::categoriseFromWoE.Wide(
6       df = loansTest,
7       varName = bestBins$variable[index],
8       woeTable = bestBins$WoE[[index]]
9     )

```

```
10 predictionCategories <- cbind(predictionCategories, binned)
11 }
12 }
```

```
1 # Retain only the relevant scorecard categories
2 predictionMatrix <-
3   predictionCategories[, !is.na(match(
4     names(predictionCategories),
5     str_remove_all(GLMCoefficients$variableName, "\`"))
6   ))] %>%
7   as.matrix()
```

```
1 predictionMatrix <- cbind(as.vector(rep.int(x = 1, times = dim(predictionMatrix)[1])),
2                             predictionMatrix)
```

```
1 TestLogit <- predictionMatrix %*% CoefficientsVector
2 TestLogit <-
3   tibble::enframe(TestLogit[, 1]) %>%
4   mutate(
5     p = 1 / (1 + exp(-value)),
6     oddsGood = if_else(is.infinite(p / (1 - p)), 1e10, p / (1 - p)))
7
8 predictionScorecard <- predictionMatrix %*% ( GLMScores$points %>% as.matrix() )
```

```
1 loansTest %>%
2   cbind(TestLogit) %>%
3   cbind(predictionScorecard) %>%
4   filter(predictionScorecard > 0) %>%
5   ggplot(aes(predictionScorecard)) +
6   geom_density(col = "blue", fill = "lightblue", adjust = 3)
```

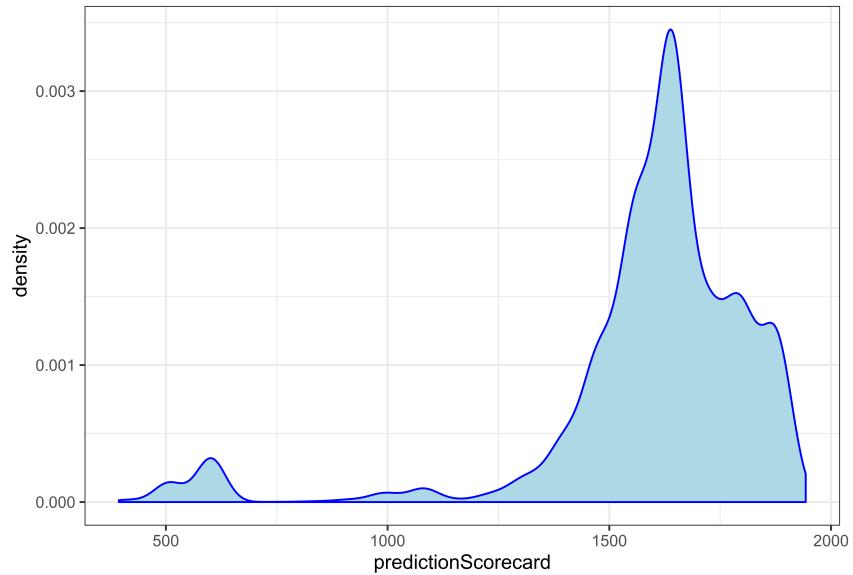


Figure 3.2: Density of loans by scorecard

We can now see the same downward dynamics as training set

```

1 loansTest %>%
2   cbind(TestLogit) %>%
3   filter(between(value, -2, 5)) %>%
4   ggplot(aes(value, col = grade)) +
5   geom_density(adjust = 2)

```

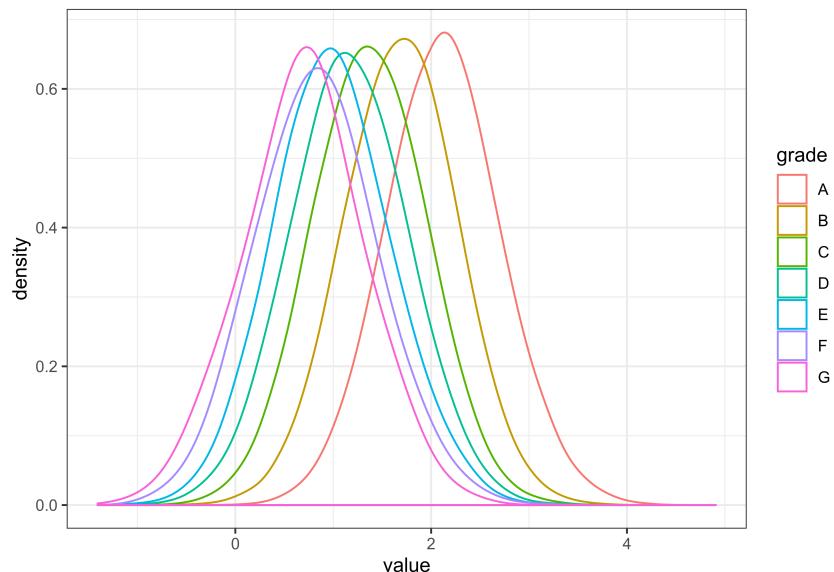


Figure 3.3: Logit value predicted by the model

```

1 loansTest %>%
2   cbind(TestLogit) %>%
3   ggplot(aes(p, col = grade)) +
4   geom_density(adjust = 2) +
5   scale_x_log10()

```

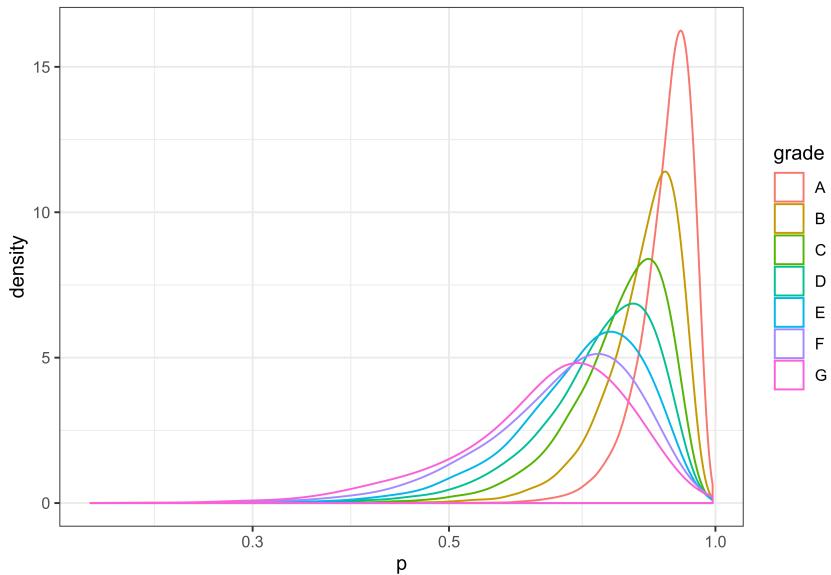


Figure 3.4: Probability of a loan being GOOD predicted by the model

```

1 loansTest %>%
2   cbind(TestLogit) %>%
3   filter(between(oddsGood, 0.1, 100)) %>%
4   ggplot(aes(oddsGood, col = grade)) +
5   geom_density(adjust = 2) +
6   scale_x_log10()

```

3.7 Confusion matrix

Given a probability p from the model, we use a $p = 0.50$ cut-off point to decide whether a loan is Good or Bad. The Confusion Matrix results are:

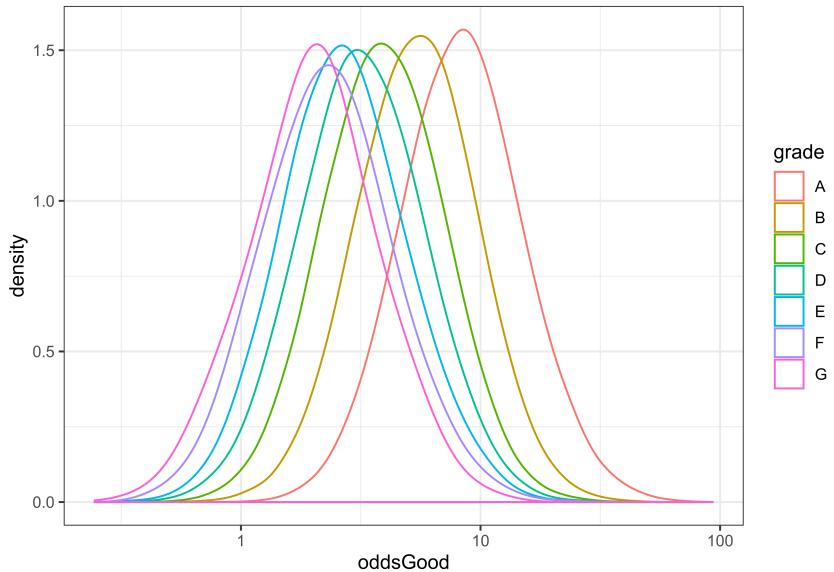


Figure 3.5: Odds of a loan being GOOD predicted by the model

```

1 tCM <- loansTest %>%
2   cbind(TestLogit) %>%
3   select(p, isGoodLoan) %>%
4
5   mutate(p = if_else(p >= 0.50, "GOOD", "BAD"),
6         isGoodLoan = if_else(isGoodLoan, "GOOD", "BAD")) %>%
7   rename(Predicted = p,
8         Actual = isGoodLoan) %>%
9
10  table() %>%
11    caret::confusionMatrix(positive = "GOOD")
12
13 tCM
14
15 ## Confusion Matrix and Statistics
16 ##
17 ##          Actual
18 ## Predicted   BAD    GOOD
19 ##      BAD     1864    1512
20 ##      GOOD   50657  207239
21 ##
22 ##          Accuracy : 0.8003
23 ##                 95% CI : (0.7988, 0.8019)
24 ##      No Information Rate : 0.799
25 ##      P-Value [Acc > NIR] : 0.043
26 ##
27 ##          Kappa : 0.0435
28 ##
29 ##  Mcnemar's Test P-Value : <2e-16
30 ##
31 ##          Sensitivity : 0.99276
32 ##          Specificity : 0.03549
33 ##          Pos Pred Value : 0.80358
34 ##          Neg Pred Value : 0.55213
35 ##          Prevalence : 0.79898
36 ##          Detection Rate : 0.79319

```

```

23 ##      Detection Prevalence : 0.98708
24 ##      Balanced Accuracy : 0.51412
25 ##
26 ##      'Positive' Class : GOOD
27 ##

```

The results suggest that the model is effective at predicting good loans (measured by the sensitivity = $\frac{TP}{TP+FN} = 99.28\%$). However, this is deceptive since the dataset is unbalanced. The model is performing poorly at detecting bad loans (measured by the specificity = $\frac{TN}{TN+FP} = 3.55\%$) The dataset is unbalanced and measuring the model's performance with a confusion matrix is imprecise. A much better approach would be to train many models on balanced datasets (by sampling a reduced 'good loans' dataset) and study the distribution of that resulting models and their parameters.

More critically, the confusion matrix does not (and cannot) reflect the consequence of getting predictions wrong. At the end of the day, the only relevant consequence is estimating the number of dollars lost on a loan. A misqualified loan might lead to a loss of a single dollar, or a million. Predicting a probability of default is not enough. We need to subsequently predict the expected loss when a particular loan defaults. In the conclusion, we suggest one possible avenue.

3.8 ROC Curve

A popular measure for the performance of such a logistic regression model is to consider its *Receiver Operating Characteristic* curve (*ROC*) and calculate its area under curve (*AUC*). We use the `ROCR` package (([Sing et al., 2005](#))). We recommend @fawcett2004roc for a very good overview of Receiver Operating Characteristics graphs.

We first create a prediction object that will

```

1 ROCRPrediction <- ROCR::prediction(TestLogit$p, loansTest$isGoodLoan)

```

Figure 3.6 plots the Receiver Operating Characteristic curve of the regression. The area under the ROC curve is 68.23%.

```

1 ROCR::performance(ROCRPrediction,
2                         measure = "tpr",
3                         x.measure = "fpr") %>%
4   ROCR::plot(colorize = TRUE)

```

The AUC has an important statistical property: the AUC of a classifier is equivalent to the probability that the classifier will rank a randomly chosen positive instance higher than a randomly

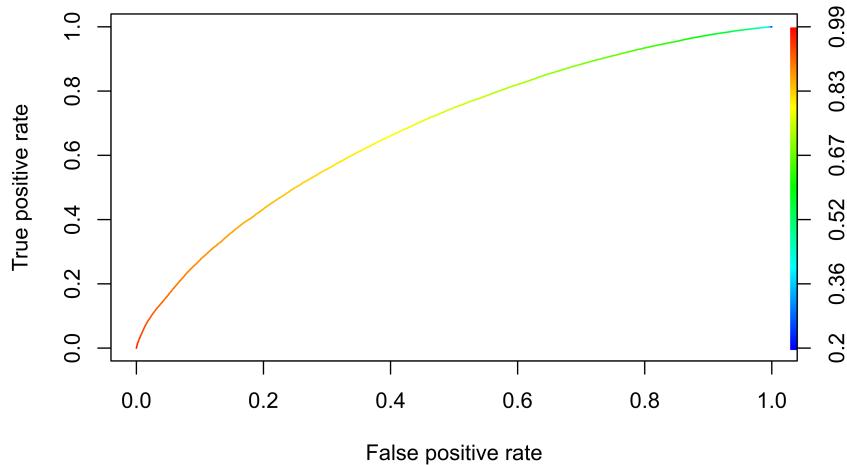


Figure 3.6: Receiver Operating Characteristic

chosen negative instance. This is equivalent to the Wilcoxon test of ranks. The AUC is also closely related to the Gini index, which is twice the area between the diagonal and the ROC curve ($\text{Gini} + 1 = 2 \times \text{AUC}$).

```

1 ROCR::performance(ROCRPrediction,
2   measure = "prec",
3   x.measure = "rec") %>%
4   ROCR::plot(colorize = TRUE)

```

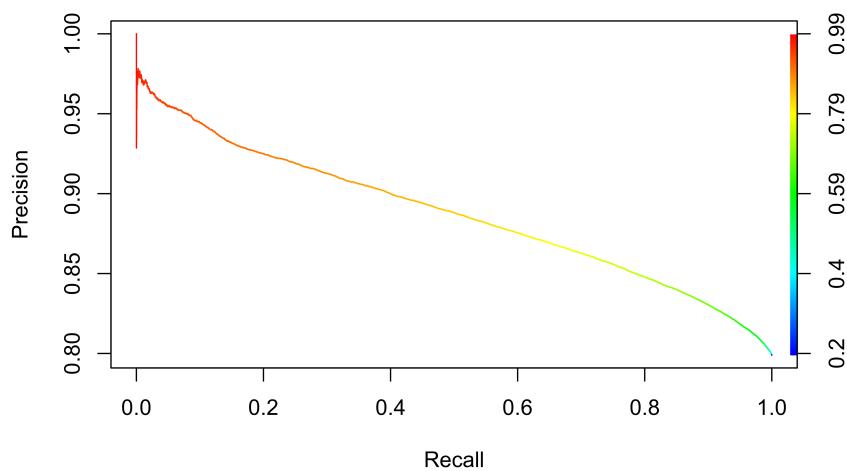


Figure 3.7: Precision/Recall curve

```

1 ROCR::performance(ROCRPrediction,
2   measure = "sens",
3   x.measure = "spec") %>%
4   ROCR::plot(colorize = TRUE)

```

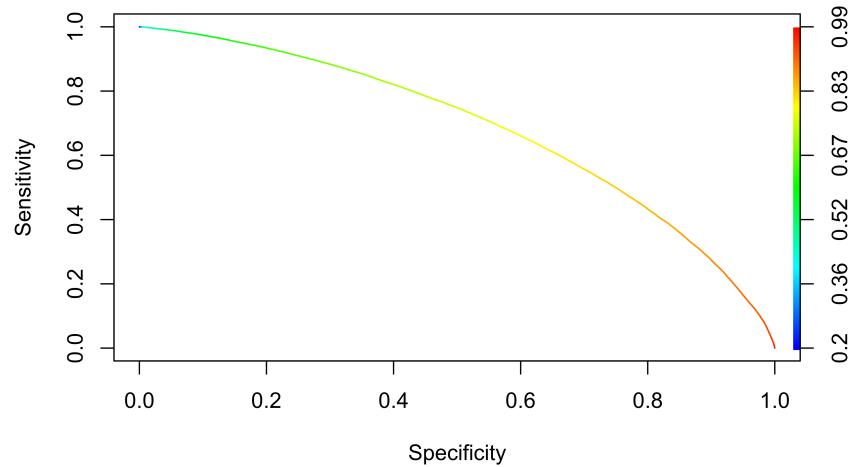


Figure 3.8: Sensitivity/Specificity curve

```

1 ROCRPerformance <- ROCR::performance(ROCRPrediction,
2   measure = "lift",
3   x.measure = "rpp") %>%
4   ROCR::plot(colorize = TRUE)

```

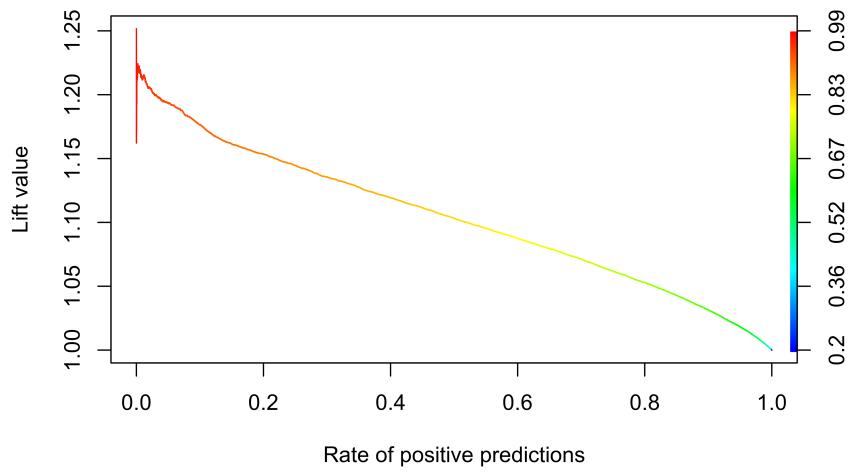


Figure 3.9: Lift Chart

Chapter 4

Conclusion

This report was an exploration of the LendingClub dataset. The result should be seen as unsatisfactory: the final model clearly extracts relevant information that classifies loan applications in a way that broadly mirrors the ratings proposed by LendingClub. It is however too imprecise to draw any comfort that it could be used in any way.

One question, which is at the end of the day the only relevant question, has not been answered: Is the interest rate proposed by LC high enough?

This report has additional clear limitations:

- Facing such a large dataset, it took a very long time before settling on a tractable question. Exploring the data lead to many blind alleys with little interest. See next section for further explanations.¹
- As a practical tool, our approach would not work in real life: we have no data on rejected loans. Using this model to accept/reject loans would need more work (for example using Reject Inference).
- We only considered probabilities of default.

Further possible avenues to explore are numerous:

- Address the unbalanced dataset by sampling a number of training samples whose size is identical to the test dataset.
- Use PCA (albeit on an extract of the dataset) to narrow the variables.
- Model the time-to-default to also provided guidance on the maximum term of the accepted loans. Maybe some sort of multivariate Poisson process (if such distributions exist).
- Extend to LGD with Good Dollars and Bad Dollars instead of Good Loans / Bad Loans. This approach is suggested in a report by SAS ((Miner, 2012)).
- Improve the regularisation of the model parameters.
- Explore economic cycles/situations as additional entry.
- We use the entire dataset to estimate the impact of time as a polynomial curve. To assess future loans, this would not be acceptable. Only an online algorithm should be used. For example ARMA/ARIMA, Kalman filtering of the time-trend trajectory?

¹But there is no Goldilocks dataset: they are always to small or too big.

- The size of the dataset is an issue to apply other techniques. But with stochastic methods, possible other models could be:
 - Tree models which have numerous variations (CART generally, and simple or aggregated boosted decision trees specifically)
 - Neural network

Chapter 5

Errands and post-mortem

This project took much longer than the MovieLens capstone. Here is a short list of ideas explored, pitfalls, lessons learned, and character-builders. This last section is written in the belief that knowing what does not work is as worthy as what does work.¹

5.1 Modeling

One bit of advice I have known of, agreed with, and naturally set aside, was to build early and small. In other words, deep and wide exploration of the dataset for the sake of it without modeling is wasteful of time and ideas. Modeling early, small and wrong (at least initially) is a better way to keep track of progress, a good sense of the eventual challenges, and suggest fruitful data exploration avenues. More importantly, it is in and by itself of form of exploration of the dataset. It also tends to explore the data wide-and-shallow rather than deep-but-narrow.

Having said that, I did some early exploration of models. Principal Component Analysis, naive Linear Regression, Extreme Boosting and Random Forest were toyed with. No model could be trained on the full set. I therefore came to limit the training set on a random sample of 0.1% (1 thousandth) of the initial full set.

I then went on a quest to formulate a model suitable for batch training (either simple sequential batches or stochastic) and gradient descent.

I also considered online training. However, an untested (and possibly completely wrong) intuition was that online methods are not adapted to an unbalanced dataset: the number of defaults/write-offs is low for high quality ratings. However, the dataset is evidently a time series which points to online training.

Along the way, the main unanswered question remained what to study out of the dataset: focus on a very narrow of variables? enrich it with other sources to study the impact of cyclical economic crises? determine an optimal pricing of each loan?

Let's look at a couple of those:

5.2 Geographical data

We sourced US zip and FIPS (*Federal Information Processing Standards*) codes, and macroeconomical data for possible geographical statistics. The source code for the data import and

¹Nobody gets a prize for having been incorrect. But nobody gets a prize not knowing where others have been incorrect.

reformatting is available in the GitHub repository.

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. Geofred turned out the best data source.

- Median income per household: <https://geofred.stlouisfed.org/map/?th=pubugn&cc=5&rc=false&im=fractile&sb&lng=-112.41&lat=44.31&zsm=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&zsm=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&zsm=4&sl&sv&am=Average&at=Not%20Seasonally%20Adjusted,%20Monthly,%20Percent&sti=1224&fq=Monthly&rt=county&un=lin&dt=2019-08-01>

Of key interest were indicators of economic stress with the following intuition: if a borrower had a good credit standing, traditional banks would provide cheaper access to credit. In times of financial distress (indicated by higher unemployment, lower GDP growth, income per household, ...), we should see higher volumes of loans, and/or changes in the loan application patterns.

Given the time available, that data turned out to be too difficult to use:

- it was incomplete since the reported figures were not available over the entire timespan of the LendingClub dataset;
- ZIP codes and FIPS location reference change over time, meaning that determining economic indicators for a given loan was no easy to automatise or would have require substantial time-consuming hand-made adjustment.

Lesson learned: data sourcing, cleaning is a full-time job by itself. (De Prado, 2018), as one of many sources, warned us. And was ignored...

5.3 Stochastic Gradient Descent

5.3.1 Conclusions

This subsection is a bit detailed. The main final points are however simple:

- A more thorough search for existing literature and craft would have been time better spent.
- Throwing more *iron* should be a last resort after trying to be more clever (or more realistically looking for what actually clever people have done in the past). This report is not intended to be a PhD thesis.
- R is not (yet) equipped to take advantage of advances in full program automatic differentiation (and leverage deep learning optimisation library) like Swift ² and Julia ³ are. R has the **Madness** package and can also interface with Julia. But that was another rabbit hole that was too dark from the outset.

²https://blog.tensorflow.org/2019/06/fastais-deep-learning-from-foundations_28.html

³<https://fluxml.ai/Zygote.jl/latest/>

- On the positive side, calculating the derivative of the multimodal NPV function was checked with Maxima symbolic math capabilities and was an opportunity to re-acquaint ourselves with it.

Here is what was written at the time of exploring SGD.

The early exploration of stochastic gradient as the only way to tackle extremely large dataset was decided on the basis of this diagram: 5.1⁴. It is otherwise extremely valuable, but we should realised much earlier that financial institutions have beeen dealing with such datasets for decades, at times when computing capacity was a orders of magnitude lower than now. If they could do it then, exploring how they did it should have been done first.

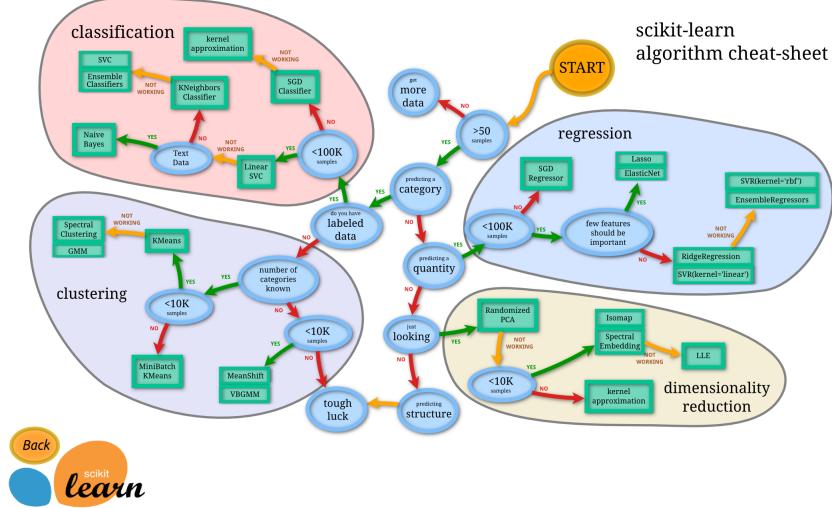


Figure 5.1: Scikit Learn algorithm cheat-sheet

5.3.2 Gradient descent

Gradient descent is a generic numerical optimisation algorithm to iteratively converge towards a (sometimes local) minimum of a given function. It is extensively used in statistical learning to minimise error functions.

In the case of a simple linear regression model, the model training error J (the *cost function*) as a function of θ is:

$$J_{(\theta)} = \frac{1}{2} \sum_{i=1}^N \epsilon(y_i, \theta X_i)$$

where the model parameters are denoted θ_i , $X_i \in \mathbb{R}^k$ are the predictors, $Y_i \in \mathbb{R}^k$ are the responses and ϵ is a distance function. Typically, ϵ will be the Manhattan error or the Euclidian norm ($A = (a_1, \dots, a_n), B = (b_1, \dots, b_n)$).

Manhattan: $\epsilon(A, B) = \sum_{i=1}^n |a_i - b_i|$

Euclidian norm: $\epsilon(A, B) = \sqrt{\sum_{i=1}^n (a_i - b_i)^2}$

The gradient descent algorithm uses the gradient of the error function, $\nabla J_{(\theta)}$, defined as:

⁴Source: https://scikit-learn.org/stable/tutorial/machine_learning_map/index.html

$$\nabla J_{(\theta)} = \left(\frac{\partial J}{\partial \theta_0}, \frac{\partial J}{\partial \theta_1}, \dots, \frac{\partial J}{\partial \theta_p} \right)$$

And in the case of linear regression is in a matrix form that can be computed efficiently:

$$\nabla J_{(\theta)} = (y^T - \theta X^T) X$$

The gradient decent algorithm finds parameters in the following manner iterating over the training samples:

$$\text{While } ||\alpha \nabla J_{(\theta)}|| > \eta, \theta := \theta - \alpha \nabla J_{(\theta)}$$

In practice, the cost function will add a penalty term to regularise the model parameters (see below).

5.3.3 Stochastic Gradient Descent

With realistic datasets, gradient descent can experience slow convergence because (1) each iteration requires calculation of the gradient for every single training example, and (2) since each individual sample is potentially very different from another, the calculated gradient may not be optimal. In such case, the gradient descent can be done using a batch of several training samples and use the average of the cost function (**batch gradient descent**). This addresses those two sources of inefficiency.

This method however still requires iterating over the entire dataset. We can instead iterate over batches of random training samples drawn from the entire dataset, instead of being drawn sequentially. This is the *stochastic gradient descent*.

Aside from the choice of the initial choice of samples, and the averaging of the cost function, the update of θ remains identical.

5.3.4 Cost function for multi-modal NPV

Cost function as a function of the Q parameter (equivalent to scorecard).

Looking back at the distribution of the NPV between the -1 and about 1.5, it is multimodal and looks like the sum of 4 log-normal distributions with modes centered on about

$$PDF(x) = \frac{1}{x\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2} \left(\frac{\log(x) - \mu}{\sigma}\right)^2\right)$$

$$\text{mode} = m = e^{\mu - \sigma^2}, \text{ therefore: } \mu = \log(\text{mode}) + \sigma^2$$

$$PDF(x) = \sqrt{\frac{e}{2\pi}} \frac{1}{x\sigma} \exp\left(-\frac{1}{2} \left(\frac{\log(\frac{x}{m}) - \sigma^2}{\sigma}\right)^2\right)$$

We will center the distribution on the mode, therefore:

$$PDF(x) = \sqrt{\frac{e}{2\pi}} \frac{1}{(x - m)\sigma} \exp\left(-\frac{1}{2} \left(\frac{\log(\frac{x-m}{m}) - \sigma^2}{\sigma}\right)^2\right)$$

The distribution's tail is towards positive infinity. For the symmetric result, we would replace $(x - m)$ by $-(x - m)$.

If we use 4 log-normal distributions, the cost function is:

$$J(x, Q) = -[x - (\alpha_1 \text{PDF}_1(x, Q) + \alpha_2 \text{PDF}_2(x, Q) + \alpha_3 \text{PDF}_3(x, Q) + \alpha_4 \text{PDF}_4(x, Q))]^2$$

To optimise the shape of the total multi-modal distribution, we will assume that each α , m and σ is a linear function of Q . The derivative $\frac{\partial J}{\partial Q}$ is⁵:

$$\begin{aligned} \frac{\partial J}{\partial Q}(x, Q) = & -\sqrt{\frac{2}{\pi}}x - \sqrt{\frac{2}{\pi}\frac{\alpha_1}{\sigma_1(-x+m_1)}}e^{-\frac{1}{2}\left(\frac{\log\left(\frac{-x+m_1}{m_1}\right)+\sigma_1^2}{\sigma_1}\right)^2} \\ & - \frac{1}{\sqrt{2\pi}}\frac{\alpha_2}{\sigma_2(-x+m_2)}e^{-\frac{1}{2}\left(\frac{\log\left(\frac{-x+m_2}{m_2}\right)+\sigma_2^2}{\sigma_2}\right)^2} \\ & - \sqrt{\frac{2}{\pi}\frac{\alpha_3}{\sigma_3(-x+m_3)}}e^{-\frac{1}{2}\left(\frac{\log\left(\frac{-x+m_3}{m_3}\right)+\sigma_3^2}{\sigma_3}\right)^2} \\ & - \sqrt{\frac{2}{\pi}\frac{\alpha_4}{\sigma_4(x-m_4)}}e^{-\frac{1}{2}\left(\frac{\log\left(\frac{x-m_4}{m_4}\right)+\sigma_4^2}{\sigma_4}\right)^2} \end{aligned}$$

We then worked on the basis of code developed for the *Movielens* capston⁶.

5.4 Final final Conclusion

We only received 999 cuts. We survived...

⁵This was actually generated using Maxima (code in Appendix) which allows for quicker iterations.

⁶See <https://github.com/Emmanuel-R8>/<https://github.com/Emmanuel-R8/HarvardX-Movielens>

Chapter 6

Appendix

6.1 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.

- *Dataset quality:* Aside from cents rounding issues, the dataset does not contain any flagrant errors that we could see (e.g. minor error of amount or rate, zipcode). Quality of the variable description is a different matter altogether.
- *Ratings:* The day-1 rating is between A1 and (and no lower than) G5. No note is rated lower than E5 after 6 November 2017, and lower than D5 after 30 June 2019.
- *Credit history:* 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.
- *Recoveries:* Recoveries (if any) are assumed to be paid 3 months after the last scheduled payment date (variable `last_pymnt_d`)
- *Survival effect:* The dataset does not include applications that were rejected by the lender (for whatever reason) or by the borrower (for example because the interest rate quote is too high). It may also be the case that some actual loans were excluded as and when the dataset changed over the years.
- LIBOR funding rate: we use the 3-year and 5-year swap rates. In reality, we should have used average tenor-weighted swap rates (i.e. ca. 1.5 Y and 2.5 Y). This requires a full swap curve and more calculation than necessary for our purpose. The principles of this report should not be significantly affected by this approximation.

We expect that LendingClub investors receive information of much better quality since the largest investors are now banks.

6.2 Data preparation and formatting

We used different sources of information:

- The LendingClub dataset made available on Kaggle;
- US geographical 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. It was automatically formatted by *RStudio*.

6.2.1 LendingClub dataset

See `Scripts/R-Data-Preparation.R`.

6.2.2 Zip codes and FIPS codes

The R package `zipcode` was installed.

```

1  #
2  # ZIPCodes dataset.
3  #
4
5  library(zipcode)
6  data(zipcode)
7  zips <- zipcode %>%
8    as_tibble() %>%
9    mutate(zip = as.integer(str_sub(zip, 1, 3)))
10
11 saveRDS(zips, "datasets/zips.rds")

```

6.2.3 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>

See `Scripts/R-LIBOR-Rates.R`.

6.3 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 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle

Variable Name	Used in model?	Description
loanID	YES	NOTE THIS IS NOT AN ORIGINAL VARIABLE. IT WAS ADDED FOR THE PURPOSE OF TRACKING LOANS INDIVIDUALLY AS AND WHEN NEEDED.
loan_amnt	YES	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	NO	The total amount committed to that loan at that point in time.
funded_amnt_inv	NO	The total amount committed by investors for that loan at that point in time.
term	YES	The number of payments on the loan. Values are in months and can be either 36 or 60.
int_rate	YES	Interest Rate on the loan
installment	NO	The monthly payment owed by the borrower if the loan originates.
grade	YES	LC assigned loan grade
sub_grade	YES	LC assigned loan subgrade
emp_title	NO	The job title supplied by the Borrower when applying for the loan.
emp_length	YES	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	YES	The home ownership status provided by the borrower during registration or obtained from the credit report. Our values are: RENT, OWN, MORTGAGE, OTHER, NONE
annual_inc	NO	The self-reported annual income provided by the borrower during registration. NOT USED AS A VARIABLE SINCE JOINT INCOME ALREADY INCLUDES IT.
verification_status	YES	Indicates if income was verified by LC, not verified, or if the income source was verified
issue_d	YES	The month which the loan was funded
loan_status	NO	Current status of the loan
pymnt_plan	NO	Indicates if a payment plan has been put in place for the loan
url	NO	URL for the LC page with listing data.
desc	NO	Loan description provided by the borrower
purpose	YES	A category provided by the borrower for the loan request.
title	NO	The loan title provided by the borrower
zip_code	NO	The first 3 numbers of the zip code provided by the borrower in the loan application.
addr_state	YES	The state provided by the borrower in the loan application

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
dti	YES	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	YES	The number of 30+ days past-due incidences of delinquency in the borrower s credit file for the past 2 years
earliest_cr_line	YES	The month the borrower s earliest reported credit line was opened
inq_last_6mths	YES	The number of inquiries in past 6 months (excluding auto and mortgage inquiries)
mths_since_last_delinq	YES	The number of months since the borrower s last delinquency.
mths_since_last_record	YES	The number of months since the last public record.
open_acc	YES	The number of open credit lines in the borrower s credit file.
pub_rec	YES	Number of derogatory public records
revol_bal	YES	Total credit revolving balance
revol_util	YES	Revolving line utilization rate, or the amount of credit the borrower is using relative to all available revolving credit.
total_acc	YES	The total number of credit lines currently in the borrower s credit file
initial_list_status	NO	The initial listing status of the loan. Possible values are – W, F
out_prncp	NO	Remaining outstanding principal for total amount funded. NOTE ONCE A LOAN IS REPAYED OR CHARGED OFF, THIS AMOUNT BECOMES 0.
out_prncp_inv	NO	Remaining outstanding principal for portion of total amount funded by investors. NOTE ONCE A LOAN IS REPAYED OR CHARGED OFF, THIS AMOUNT BECOMES 0.
total_pymnt	NO	Payments received to date for total amount funded
total_pymnt_inv	NO	Payments received to date for portion of total amount funded by investors
total_rec_prncp	NO	Principal received to date. NOTE THIS AMOUNT WILL SHOW WHETHER A BORROWER DID NOT REPAY IN FULL
total_rec_int	NO	Interest received to date
total_rec_late_fee	NO	Late fees received to date
recoveries	NO	Post charge off gross recovery. NOTE IF A LOAN IS REPAYED, THIS AMOUNT IS 0.
collection_recovery_fee	NO	Post charge off collection fee
last_pymnt_d	NO	Last month payment was received
last_pymnt_amnt	NO	Last total payment amount received
next_pymnt_d	NO	Next scheduled payment date

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
last_credit_pull_d	NO	The most recent month LC pulled credit for this loan
collections_12_mths_ex_med	NO	Number of collections in 12 months excluding medical collections
mths_since_last_major_derog	NO	Months since most recent 90-day or worse rating
policy_code	NO	Publicly available policy_code=1 / New products not publicly available policy_code=2
application_type	YES	Indicates whether the loan is an individual application or a joint application with two coborrowers
annual_inc_joint	YES	The combined self-reported annual income provided by the coborrowers during registration
dti_joint	YES	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	YES	Indicates if income was verified by LC, not verified, or if the income source was verified
acc_now_delinq	YES	The number of accounts on which the borrower is now delinquent.
tot_coll_amt	NO	Total collection amounts ever owed
tot_cur_bal	NO	Total current balance of all accounts
open_acc_6m	NO	Number of open trades in last 6 months
open_act_il	NO	Number of currently active installment trades
open_il_12m	NO	Number of installment accounts opened in past 12 months
open_il_24m	NO	Number of installment accounts opened in past 24 months
mths_since_rcnt_il	NO	Months since most recent instalment accounts opened
total_bal_il	NO	Total current balance of all installment accounts
il_util	NO	Ratio of total current balance to high credit/credit limit on all install acct
open_rv_12m	YES	Number of revolving trades opened in past 12 months
open_rv_24m	YES	Number of revolving trades opened in past 24 months
max_bal_bc	YES	Maximum current balance owed on all revolving accounts
all_util	NO	Balance to credit limit on all trades
total_rev_hi_lim	NO	Total revolving high credit/credit limit
inq_fi	YES	Number of personal finance inquiries
total_cu_tl	NO	Number of finance trades
inq_last_12m	NO	Number of credit inquiries in past 12 months
acc_open_past_24mths	NO	Number of trades opened in past 24 months.
avg_cur_bal	YES	Average current balance of all accounts

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
bc_open_to_buy	YES	Total open to buy on revolving bankcards.
bc_util	YES	Ratio of total current balance to high credit/credit limit for all bankcard accounts.
chargeoff_within_12_mths	NO	Number of charge-offs within 12 months
delinq_amnt	NO	The past-due amount owed for the accounts on which the borrower is now delinquent.
mo_sin_old_il_acct	YES	Months since oldest bank instalment account opened
mo_sin_old_rev_tl_op	YES	Months since oldest revolving account opened
mo_sin_rcnt_rev_tl_op	YES	Months since most recent revolving account opened
mo_sin_rcnt_tl	YES	Months since most recent account opened
mort_acc	YES	Number of mortgage accounts.
mths_since_recent_bc	YES	Months since most recent bankcard account opened.
mths_since_recent_bc_dlq	YES	Months since most recent bankcard delinquency
mths_since_recent_inq	YES	Months since most recent inquiry.
mths_since_recent_revol_delinq	YES	Months since most recent revolving delinquency.
num_accts_ever_120_pd	YES	Number of accounts ever 120 or more days past due
num_actv_bc_tl	YES	Number of currently active bankcard accounts
num_actv_rev_tl	YES	Number of currently active revolving trades
num_bc_sats	YES	Number of satisfactory bankcard accounts
num_bc_tl	YES	Number of bankcard accounts
num_il_tl	YES	Number of installment accounts
num_op_rev_tl	YES	Number of open revolving accounts
num_rev_accts	YES	Number of revolving accounts
num_rev_tl_bal_gt_0	YES	Number of revolving trades with balance >0
num_sats	YES	Number of satisfactory accounts
num_tl_120dpd_2m	YES	Number of accounts currently 120 days past due (updated in past 2 months)
num_tl_30dpd	YES	Number of accounts currently 30 days past due (updated in past 2 months)
num_tl_90g_dpd_24m	YES	Number of accounts 90 or more days past due in last 24 months
num_tl_op_past_12m	YES	Number of accounts opened in past 12 months
pct_tl_nvr_dlq	YES	Percent of trades never delinquent
percent_bc_gt_75	YES	Percentage of all bankcard accounts > 75% of limit.
pub_rec_bankruptcies	YES	Number of public record bankruptcies
tax_liens	YES	Number of tax liens
tot_hi_cred_lim	YES	Total high credit/credit limit
total_bal_ex_mort	YES	Total credit balance excluding mortgage
total_bc_limit	YES	Total bankcard high credit/credit limit
total_il_high_credit_limit	YES	Total installment high credit/credit limit
revol_bal_joint	YES	Total credit revolving balance

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
sec_app_earliest_cr_line	NO	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	NO	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	NO	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	NO	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	NO	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	NO	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	NO	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_no_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	NO	Flags whether or not the borrower is on a hardship plan
hardship_type	NO	Describes the hardship plan offering

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
hardship_reason	NO	Describes the reason the hardship plan was offered
hardship_status	NO	Describes if the hardship plan is active, pending, cancelled, completed, or broken
deferral_term	NO	Amount of months that the borrower is expected to pay less than the contractual monthly payment amount due to a hardship plan
hardship_amount	NO	The interest payment that the borrower has committed to make each month while they are on a hardship plan
hardship_start_date	NO	The start date of the hardship plan period
hardship_end_date	NO	The end date of the hardship plan period
payment_plan_start_date	NO	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	NO	The number of months the borrower will make smaller payments than normally obligated due to a hardship plan
hardship_dpd	NO	Account days past due as of the hardship plan start date
hardship_loan_status	NO	Loan Status as of the hardship plan start date
orig_projected_additional_amnt	NO	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_amnt	NO	The payoff balance amount as of the hardship plan start date
hardship_last_payment_amnt	NO	The last payment amount as of the hardship plan start date
disbursement_method	YES	The method by which the borrower receives their loan. Possible values are: CASH, DIRECT_PAY
debt_settlement_flag	NO	Flags whether or not the borrower, who has charged-off, is working with a debt-settlement company.
debt_settlement_flag_date	NO	The most recent date that the Debt_Settlement_Flag has been set
settlement_status	NO	The status of the borrower's settlement plan. Possible values are: COMPLETE, ACTIVE, BROKEN, CANCELLED, DENIED, DRAFT
settlement_date	NO	The date that the borrower agrees to the settlement plan
settlement_amount	NO	The loan amount that the borrower has agreed to settle for
settlement_percentage	NO	The settlement amount as a percentage of the payoff balance amount on the loan

Table 6.1: Description of the dataset variables as provided in the dataset downloaded from Kaggle (*continued*)

Variable Name	Used in model?	Description
settlement_term	NO	The number of months that the borrower will be on the settlement plan

6.4 Maxima derivation of the cost function

```

PDF1(x, Q) := alpha1( Q) * sqrt( 1 / ( 2 * pi)) *
              exp( - 1 / 2*(( log( -( x - m1( Q)) / m1( Q)) + sigma1( Q) ^ 2) /
              sigma1( Q)) ^ 2) / ( -( x - m1(Q)) * sigma1( Q)) ;

PDF2(x, Q) := alpha2( Q) * sqrt( 1 / ( 2 * pi)) *
              exp( - 1 / 2*(( log( -( x - m2( Q)) / m2( Q)) + sigma2( Q) ^ 2) /
              sigma2( Q)) ^ 2) / ( -( x - m2( Q)) * sigma2( Q)) ;

PDF3(x, Q) := alpha3( Q) * sqrt( 1 / ( 2 (* pi)) *
              exp( - 1 / 2*(( log( -( x - m3( Q)) / m3( Q)) + sigma3( Q) ^ 2) /
              sigma3( Q)) ^ 2) / ( -( x - m3( Q)) * sigma3( Q)) ;

PDF4(x, Q) := alpha4( Q) * sqrt( 1 / ( 2 * pi)) *
              exp( - 1 / 2*(( log( ( x - m4( Q)) / m4( Q)) + sigma4( Q) ^ 2) /
              sigma4( Q)) ^ 2) / ( ( x - m4( Q)) * sigma4( Q)) ;

alpha1(Q) := am1* Q + an1 ;
alpha2(Q) := am2* Q + an2 ;
alpha3(Q) := am3* Q + an3 ;
alpha4(Q) := am4* Q + an4 ;

m1(Q) := mm1* Q + mn1 ;
m2(Q) := mm2* Q + mn2 ;
m3(Q) := mm3* Q + mn3 ;
m4(Q) := mm4* Q + mn4 ;

sigma1(Q) := sm1* Q + sn1 ;
sigma2(Q) := sm2* Q + sn2 ;
sigma3(Q) := sm3* Q + sn3 ;
sigma4(Q) := sm4* Q + sn4 ;

J(x, Q): = -( x - ( PDF1( x, Q) + PDF2( x, Q) + PDF3( x, Q) + PDF4( x, Q))) ^ 2 ;

diff( PDF1(x, Q), Q) ;

```

6.5 System version

name	value
sysname	Linux
release	5.3.0-24-generic
version	#26-Ubuntu SMP Thu Nov 14 01:33:18 UTC 2019
nodename	x260
machine	x86_64

	Loaded Package
1	binner
2	bookdown
3	knitr
4	dslabs
5	kableExtra
6	gridExtra
7	lubridate
8	forcats
9	stringr
10	dplyr
11	purrr
12	readr
13	tidyverse
14	tibble
15	ggplot2
16	tidyverse
17	stats
18	graphics
19	grDevices
20	utils
21	datasets
22	methods
23	base

Bibliography

- Bokhari, M. M. (2019). Credit risk analysis in peer to peer lending data set: Lending club.
- 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].
- De Prado, M. L. (2018). *Advances in financial machine learning*. John Wiley & Sons.
- Hothorn, T., Hornik, K., and Zeileis, A. (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3):651–674.
- Hothorn, T. and Zeileis, A. (2015). partykit: A modular toolkit for recursive partytioning in r. *Journal of Machine Learning Research*, 16(118):3905–3909.
- Kan, L. C. P. W. (2019). Kaggle - LendingClub dataset. <https://www.kaggle.com/wendykan/lending-club-loan-data>. [Note: Accessed 31 August 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.
- Miner, S. E. (2012). Sas global forum 2012 financial services paper 141-2012 building loss given default scorecard using weight of evidence bins in sas®. <https://support.sas.com/resources/papers/proceedings12/141-2012.pdf>.
- Peng, R. (2012). *Exploratory data analysis with R*. Lulu. com.
- Sing, T., Sander, O., Beerenwinkel, N., and Lengauer, T. (2005). ROCR: visualizing classifier performance in R. *Bioinformatics*, 21(20):3940–3941.