



QuantUniversity, LLC

www.quantuniversity.com



Machine Learning applications for Credit Risk

Presented By:

Sri Krishnamurthy, CFA, CAP

www.QuantUniversity.com

sri@quantuniversity.com

Credit Risk and scoring



Credit risk in consumer credit

Credit-scoring models and techniques assess the risk in lending to customers.

Typical decisions:

- Grant credit/not to new applicants
- Increasing/Decreasing spending limits
- Increasing/Decreasing lending rates
- What new products can be given to existing applicants ?



Working with mixed-data

- Gower similarity coefficient is a composite measure
- It takes quantitative (such as rating scale), binary (such as present/absent) and nominal (such as worker/teacher/clerk) variables.
- Gower distance is used for calculating distances when we have mixed types of variables (continuous and categorical)
- A linear combination using user-specified weights (most simple an average) is calculated to create the final distance matrix.
- The metrics used for each data type are described below:
 - Quantitative: range-normalized Manhattan distance
 - Ordinal: variable is first ranked, then Manhattan distance is used with a special adjustment for ties
 - Nominal: variables of k categories are first converted into k binary columns and then the Dice coefficient is used
(https://en.wikipedia.org/wiki/S%C3%B8rensen%E2%80%93Dice_coefficient)



Support in R

- daisy
- Pam
- agnes

<https://stat.ethz.ch/R-manual/R-devel/library/cluster/html/pam.html>

[https://en.wikibooks.org/wiki/Data Mining Algorithms In R/Clustering/CLARA](https://en.wikibooks.org/wiki/Data_Mining_Algorithms_In_R/Clustering/CLARA)



Getting to know the Lending Club Dataset



Lending club


[Personal Loans](#)
[Business Loans](#)
[Patient Solutions](#)
[Investing](#)
[How It Works](#)
[About Us](#)
[Sign in](#)
[Help](#)

Better Rates. Together.



Personal Loans up to \$40,000



Respond to mail

Check Your Rate

Won't impact your credit score



Privacy & security
PROTECTION

Financial Innovation

Lending Club is the world's largest online marketplace connecting borrowers and investors. We're transforming the banking system to make credit more affordable and investing more rewarding. We operate at a lower cost than traditional bank lending programs and pass the savings on to borrowers in the form of lower rates and to investors in the form of solid returns.



Featured Borrower



Rebecca
Durham, NC
[Pay Off Credit Cards](#)
\$5,000 at 9.98% APR

"The application process was simple, and it was neat to watch the process as people invested in my loan."



The Data



[Personal Loans](#) |
 [Business Loans](#) |
 [Patient Solutions](#) |
 [Investing](#) |
 [How It Works](#) |
 [About Us](#)

[Sign in](#)
[Help](#)

Lending Club Statistics



Platform: [Highlights](#) |
 Public Note Offering: [Investor Performance](#) |
 [Loan Statistics](#) |
 [Download Data](#)

Want to slice and dice the data? Help yourself to the following exports of our loan databases.

DOWNLOAD LOAN DATA

Year:
 Format:

[Download](#)

These files contain complete loan data for all loans issued through the time period stated, including the current loan status (Current, Late, Fully Paid, etc.) and latest payment information. The file containing loan data through the "present" contains complete loan data for all loans issued through the previous completed calendar quarter. [Sign in](#) to download the full version of the files.

DECLINED LOAN DATA

Year:
 Format:

[Download](#)

These files contain the list and details of all loan applications that did not meet Lending Club's credit underwriting policy.

DATA DICTIONARY

The Data Dictionary includes definitions for all the data attributes included in the Historical data file and the In Funding data file.

[Format XLSX \(22kb\)](#)

<https://www.lendingclub.com/info/download-data.action>



The Data

The screenshot shows the Kaggle dataset page for 'Lending Club Loan Data'. The page header includes the Kaggle logo and navigation links for Competitions, Datasets, Kernels, Forums, and Jobs. There are 'Sign Up' and 'Log In' buttons. The dataset title 'Lending Club Loan Data' is prominently displayed, along with the subtitle 'Analyze Lending Club's issued loans' and the Lending Club logo. Below the title, it says 'by Wendy Kan · last updated 5 months ago'. A tab bar shows 'Overview' (selected), 'Kernels', 'Discussion', 'Activity', 'Download (240 MB)', 'New Notebook', and 'New Script'. The 'Overview' tab is active, showing three sections: 'Kernels', 'Discussion', and 'Top Contributors'. The 'Kernels' section lists three kernels: 'Initial loan book analysis' (56 votes, run 3 months ago), 'Python for Padawans' (15 votes, run 3 months ago), and 'Total loans by state barplot' (8 votes, run 4 months ago). The 'Discussion' section lists three discussions: 'What hardware are you using?' (1 reply, 5 days ago), 'Initial loan book analysis' (22 replies, 5 days ago), and 'Lending club model' (0 replies, 6 days ago). The 'Top Contributors' section lists three contributors: Eryk Walczak (1st), Evan Miller (2nd), and Spider Pig (3rd). At the bottom, there is a 'Recent Activity' section showing five recent events: 'adarsh2807 created kernel Notebookd2c7c92b75 28 minutes ago', 'VikramadityaArora forked kernel Python for Padawans 4 days ago', 'Abhijay Arora commented on dataset discussion What hardware are you using? 5 days ago', 'Silvio commented on kernel Initial loan book analysis 5 days ago', and 'HuangRui ran version 2 of kernel Lending Club Data insights 71a95c02 6 days ago'.



Variable description

Variable name	Variable type	Descriptions
addr_state	factor	Name of states
annual_inc	numerical	The self-reported annual income provided by the borrower during registration.
delinq_2yrs	numerical	The number of 30+ days past-due incidences of delinquency in the borrower's credit file for the past 2 years
dti	numerical	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.
emp_length	factor	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.
funded_amnt	numerical	The total amount committed to that loan at that point in time.
funded_amnt_inv	numerical	The total amount committed by investors for that loan at that point in time.
grade	factor	LC assigned loan grade
home_ownership	factor	The home ownership status provided by the borrower during registration. Our values are: RENT, OWN, MORTGAGE, OTHER.
inq_last_6mths	numerical	The number of inquiries in past 6 months (excluding auto and mortgage inquiries)
installment	numerical	The monthly payment owed by the borrower if the loan originates.
int_rate	numerical	Interest Rate on the loan
issue_d	factor	The month which the loan was funded
loan_amnt	numerical	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.
loan_status	factor	Current status of the loan
purpose	factor	A category provided by the borrower for the loan request.
sub_grade	factor	LC assigned loan subgrade
term	factor	The number of payments on the loan. Values are in months and can be either 36 or 60.
verification_status	factor	Current verification status



Objective

- Calculate dissimilarity between observations.
- Select algorithm to group observations together
- Choose the best number of clusters
- Visualize clusters on reduced dimensions



Selecting number of clusters

- Partitioning around medoids (PAM) is used in this case.
- PAM is an iterative clustering procedure with the following steps:
 - Step 1: Choose k random entities to become the medoids.
 - Step 2: Assign every entity to its closest medoid (using the distance matrix we have calculated).
 - Step 3: For each cluster, identify the observation that would yield the lowest average distance if it were to be re-assigned as the medoid. If so, make this observation the new medoid.
 - Step 4: If at least one medoid has changes, return to step 2. Otherwise, end the algorithm.



Visualization with reduced dimension

- One way to visualize many variables in a lower dimensional space is with t-distributed stochastic neighborhood embedding (t-SNE)
- This method is a dimension reduction technique that tries to preserve local structure so as to make clusters visible in a 2D or 3D visualization.
- https://en.wikipedia.org/wiki/T-distributed_stochastic_neighbor_embedding

