

Factors of Loan Defaults and Loan Application Approvals and Rejections in the United States

Introduction

Since the 2008 Housing Crisis, lending standards have become more stringent “A few things have changed since the real estate meltdown a few years ago. For purchase transactions, real estate agents will first want to know if you can get a loan. In the old days, financial institutions were doling out money to anyone with a heartbeat. Unfortunately, soft lending standards helped fuel an eventual rash of foreclosures. Suffice it to say, conditions on the ground have changed since then. Today, the best way to approach a real estate agent is with a lender pre-approval in hand. It shows that you’re ready and able to buy.” (Mariotti, 2018). To get a pre-approval, the process is “pulling a three-bureau credit report (called a tri-merge) that shows your credit score and credit history as reported by third-party, respected institutions. Within the credit report, a lender can see your payment history (to see if payment obligations have been on-time and in-full) and your lines of credit (past and present).” (Mariotti, 2018)

The reason for this process of pulling your credit information is that loan defaults costs financial institutions money, “The statistical findings showed significantly that proper management of loans given to clients will yield more profits for the firms. Also, there was a significant relationship between the problem of recovery and overdue loans and profitability. The relationship between deficient analysis of project viability and profitability was also positive. Lastly, there was a significant relationship between the problem of recovery and overdue of loans and deficient analysis of project viability.” (Ntiamoah, 2014). This extra cost cuts into financial institutions’ bottom line. Having a robust and predictive model of defaults will save the

financial institution money and ensure that the loan applicant is not provided with a loan bigger than they can afford.

Lenders review loan applications to determine which loans get approved and which do not. To make this assessment, lenders would like to see the applicant's payment history to understand their payment behavior. However, based on this behavior, they cannot always predict which loan applicants will default on their loan payment and cost the company money.

There is an opportunity to look at past loan applicants' loan payment history and if they defaulted, and then input the new loan applicant's information to predict if they will also default based on that historical data. If the prediction is that the loan applicant will default, then the loan application will be rejected. If the prediction is that the applicant will not default, then it's approved. This will help the lender determine if the loan applicant is too risky (I.e., they will default on their loan) to provide loan application approval. This will in turn save the company money.

To best predict loan defaults, a dataset was obtained from UCI Machine Learning Repository (<https://archive.ics.uci.edu/ml/datasets/default+of+credit+card+clients#>) which has a large volume of records (30,000) and 25 columns or features, which are:

Variable	Definition	Key
ID	Unique identifier of the loan application	Continuous numerical variable (no units)
LIMIT_BAL	Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit	continuous numerical value in dollars
SEX	Loan applicant gender	1 = Male

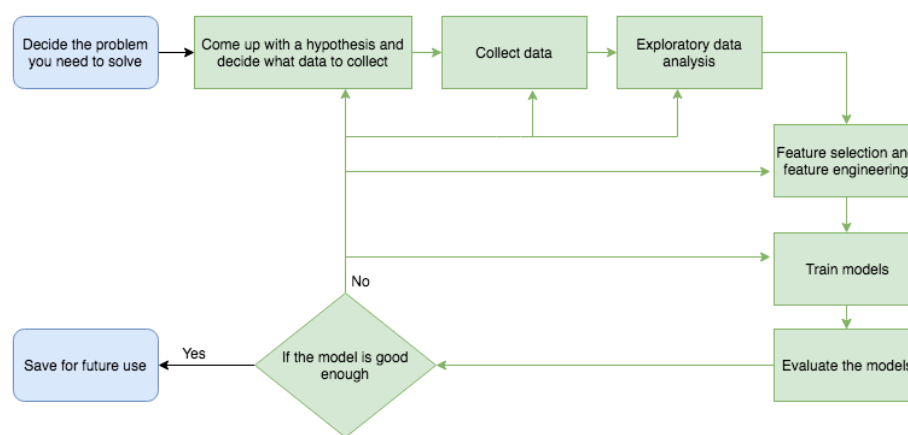
		2 = Female
EDUCATION	Loan applicant level of education	1 = graduate school; 2 = university; 3 = high school; 4 = others
MARRIAGE	Loan applicant marital status	1 = married; 2 = single; 3 = others
AGE	Loan applicant age	Continuous numerical value in years
PAY_0-6	History of past payment. We tracked the past monthly payment records (from April to September, 2005) as follows: X6 = the repayment status in September, 2005; X7 = the repayment status in August, 2005; . . .; X11 = the repayment status in April, 2005	0 = pay duly; 1 = payment delay for one month; 2 = payment delay for two months; . . .; 8 = payment delay for eight months; 9 = payment delay for nine months and above.
BILL_AMT1-6	Amount of bill statement (NT dollar). X12 = amount of bill statement in September, 2005; X13 = amount of bill statement in August, 2005; . . .; X17 = amount of bill statement in April, 2005	continuous numerical value in dollars
PAY_AMT1-6	Amount of previous payment (NT dollar). X18 = amount paid in September, 2005; X19 = amount paid in August, 2005; . . .; X23 = amount paid in April, 2005	continuous numerical value in dollars
default payment next month	If the loan applicant defaulted on the payment next month	Yes = 1, No = 0

Methods

Multiple Logistic Regression models with an 80% train and 20% test sets were used to determine which loan default dataset features are most related to or most predictive of loan defaults, where the “default payment next month” feature is the target. Because the Target Variable (default payment next month) is a binary variable, in this research initiative Logistic

Regression will be used as the model of choice. This was a supervised learning model where “the supervisor is the target variable, a column in the data representing values to predict from other columns in the data. The target variable is chosen to represent the answer to a question the organization would like to answer or a value unknown at the time the model is used that would help in decisions. Sometimes supervised learning is also called predictive modeling. The primary predictive modeling algorithms are classification for categorical target variables or regression for continuous target variables.” (Abbott, 2014, p. 5).

In addition to the logistic regression models with 1 target and 23 features, another logistic regression model was used with the 5 best features only derived with highest chi-squared statistics. See below for the process used to evaluate model success:



(Keheo, 2018)

Logistic models are not perfect and have their disadvantages, “Logistic Regression is a statistical analysis model that attempts to predict precise probabilistic outcomes based on independent features. On high dimensional datasets, this may lead to the model being over-fit on the training set, which means overstating the accuracy of predictions on the training set and

thus the model may not be able to predict accurate results on the test set. This usually happens in the case when the model is trained on little training data with lots of features. So on high dimensional datasets, Regularization techniques should be considered to avoid over-fitting (but this makes the model complex). Very high regularization factors may even lead to the model being under-fit on the training data.” (Grover, 2020). If overfitting is identified as an issue via cross-validation, two types of regularization techniques will be used: Lasso (L1 Regularization) or Ridge (L2 Regularization), where λ is called the regularization parameter and controls the trade-off between fitting the training data well and keeping the parameters small to avoid overfitting.

To measure the models’ performance, the accuracy, precision, recall and F1 score will be calculated for the logistic and decision tree models for both the 23 features and 5 best 5 features as well as visualize a confusion matrix and an ROC curve where "The area under the ROC curve (AUC) results were considered excellent for AUC values between 0.9-1, good for AUC values between 0.8-0.9, fair for AUC values between 0.7-0.8, poor for AUC values between 0.6-0.7 and failed for AUC values between 0.5-0.6." (El Khouli, 2009). All these together were used to evaluate the results of the logistic regression models. The definitions of accuracy, precision, recall, F1 and AUC scores can be found below:

- **Accuracy:** “Accuracy represents the number of correctly classified data instances over the total number of data instances” (B, H. N., 2020)
- **Precision:** “Precision talks about how precise/accurate your model is out of those predicted positive, how many of them are actual positive.” (B, H. N., 2020)
- **Recall:** “Recall actually calculates how many of the Actual Positives our model capture through labeling it as Positive (True Positive).” (B, H. N., 2020)

- **F1 Score:** “F1-score is a metric which takes into account both precision and recall.
 - **Precision:** positive value
 - **Recall:** true positive rate” (B, H. N., 2020)
- **AUC Score:** ““So, what area under the ROC curve describes good discrimination?

Unfortunately there is no "magic" number, only general guidelines. In general, we use the following rule of thumb:

- 0.5 = This suggests no discrimination, so we might as well flip a coin.
- 0.5-0.7 = We consider this poor discrimination, not much better than a coin toss.
- 0.7-0.8 = Acceptable discrimination
- 0.8-0.9= Excellent discrimination
- >0.9 = Outstanding discrimination" (Dixon, 2020)

The data was prepared and cleaned to support the creation of several models and evaluation of their performance and results for interpretation (see list below). With these results, a conclusion and recommendations in reference to the original question and problem will be detailed.

Data Preparation Steps:

- Change first row to header
- Change target variable from “default payment next month” to “DEFAULT”
- Change numerical and target variables to a numeric column type in the pandas dataframe (float64)

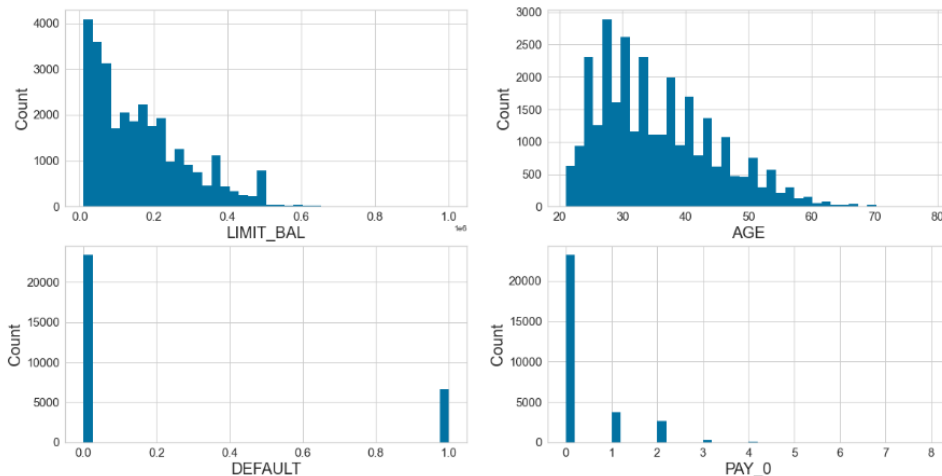
Analysis

Histograms, bar charts, Pearson's correlation matrix (numerical data) and Spearman's correlation matrix (categorical & numerical data) were created to describe the data.

- **Numerical Data**

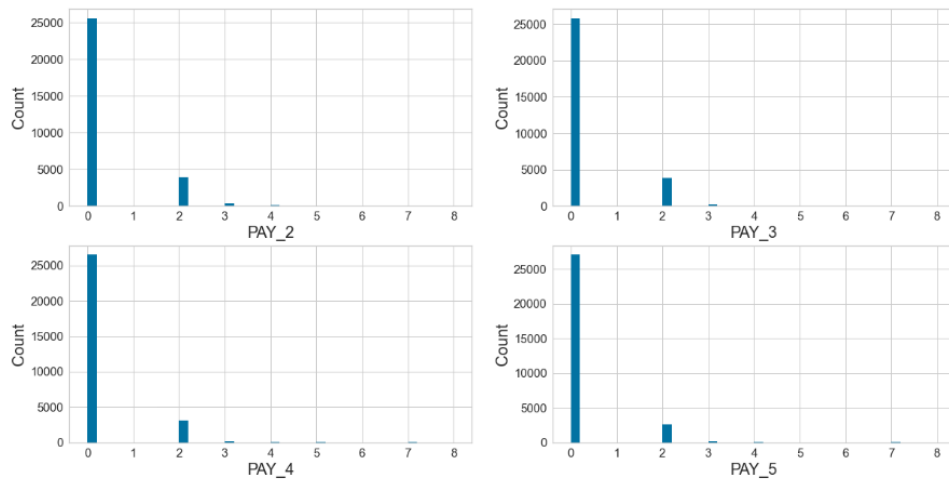
- **Histogram:** To show the distribution of the numerical data ('LIMIT_BAL', 'AGE', 'DEFAULT', 'PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5', 'BILL_AMT6', 'PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6'), histograms were used, and observations were captured:

- Most people in the dataset have a Limit Balance between 4000 and 1000 dollars
- Most people in the dataset are between the ages of 25 and 45 years
- Most people in the dataset do not default on their next loan payment
- Most people in the dataset pay their loan payments on time in the first month

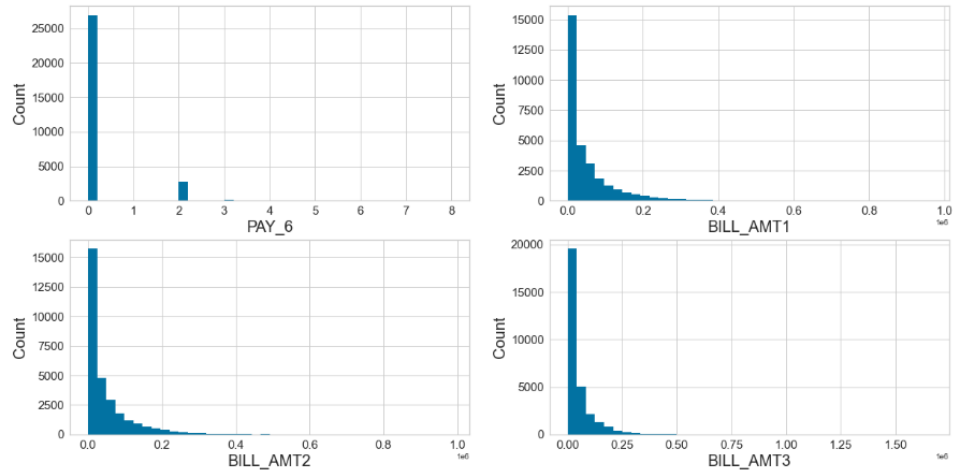


- Most people in the dataset pay their loan payments on time in the second month

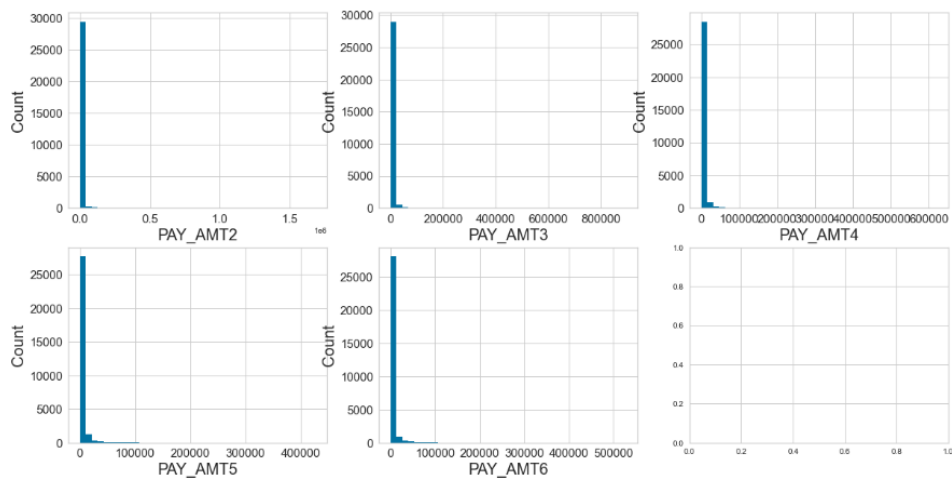
- Most people in the dataset pay their loan payments on time in the third month
- Most people in the dataset pay their loan payments on time in the fourth month
- Most people in the dataset pay their loan payments on time in the fifth month



- Most people in the dataset pay their loan payments on time in the sixth month
- Most people in the dataset have a bill payment amount between 0 and 20000 dollars in the first month
- Most people in the dataset have a bill payment amount between 0 and 20000 dollars in the second month
- Most people in the dataset have a bill payment amount between 0 and 25000 dollars in the third month



-
- Most people in the dataset have a bill payment amount between 0 and 20000 dollars in the fourth month
- Most people in the dataset have a bill payment amount between 0 and 20000 dollars in the fifth month
- Most people in the dataset have a bill payment amount between 0 and 20000 dollars in the sixth month
- Most people in the dataset have a loan payment amount between 0 and 100000 dollars in the first month



- Most people in the dataset have a loan payment amount between 0 and 25000 dollars in the second month
 - Most people in the dataset have a loan payment amount between 0 and 100000 dollars in the third month
 - Most people in the dataset have a loan payment amount between 0 and 50000 dollars in the fourth month
 - Most people in the dataset have a loan payment amount between 0 and 50000 dollars in the fifth month
 - Most people in the dataset have a loan payment amount between 0 and 50000 dollars in the sixth month
- A **Pearson's Correlation Matrix** was then created to understand the correlation between “default payment next month” and the numerical features ('LIMIT_BAL', 'AGE', 'DEFAULT', 'PAY_0', 'PAY_2', 'PAY_3', 'PAY_4', 'PAY_5', 'PAY_6', 'BILL_AMT1', 'BILL_AMT2', 'BILL_AMT3', 'BILL_AMT4', 'BILL_AMT5', 'BILL_AMT6', 'PAY_AMT1', 'PAY_AMT2', 'PAY_AMT3', 'PAY_AMT4', 'PAY_AMT5', 'PAY_AMT6') and observations were captured:
- The numerical feature most highly correlated with 'DEFAULT' is 'PAY_0', followed by PAY_2-PAY_6
 - There is not significant collinearity between the features (Rule of thumb: If the correlation > 0.8 then severe multicollinearity may be present)



- The categorical features most highly correlated with 'DEFAULT' are 'SEX_1' (male), 'EDUCATION_2' (university), 'EDUCATION_3' (high school), and 'MARRIAGE_1'(married)
- There is not significant collinearity between the features (Rule of thumb: If the correlation > 0.8 then severe multicollinearity may be present)

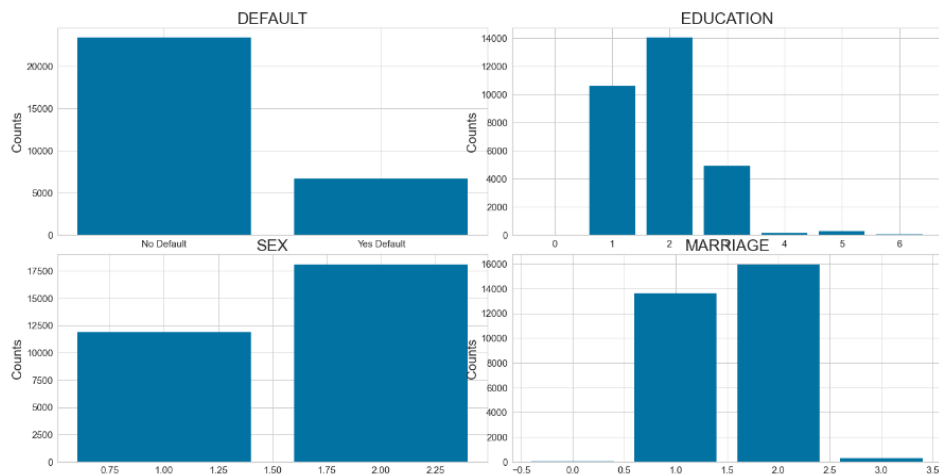
	DEFAULT	LIMIT_BAL	AGE	PAY_0	PAY_2	PAY_3	PAY_4	PAY_5	PAY_6	BILL_AMT1	BILL_AMT2	BILL_AMT3	BILL_AM
DEFAULT	1.000000	-0.169586	0.005149	0.391160	0.338456	0.294262	0.276443	0.267795	0.247234	-0.025796	-0.015970	-0.013021	-0.0090
LIMIT_BAL	-0.169586	1.000000	0.186485	-0.182430	-0.223767	-0.217512	-0.206347	-0.191925	-0.186558	0.054458	0.048928	0.061042	0.0735
AGE	0.005149	0.186485	1.000000	-0.012847	-0.020108	-0.024351	-0.015851	-0.025172	-0.029488	0.001078	0.001545	0.001952	-0.0032
PAY_0	0.391160	-0.182430	-0.012847	1.000000	0.681136	0.456119	0.397391	0.373707	0.333656	-0.071413	-0.034642	-0.022792	-0.0078
PAY_2	0.338456	-0.223767	-0.020108	0.681136	1.000000	0.634163	0.486952	0.449377	0.402223	0.076551	0.072875	0.082073	0.0929
PAY_3	0.294262	-0.217512	-0.024351	0.456119	0.634163	1.000000	0.633888	0.490391	0.440377	0.038651	0.071832	0.064948	0.0831
PAY_4	0.276443	-0.206347	-0.015851	0.397391	0.486952	0.633888	1.000000	0.669643	0.504116	0.048098	0.065707	0.089459	0.0952
PAY_5	0.267795	-0.191925	-0.025172	0.373707	0.449377	0.490391	0.669643	1.000000	0.668401	0.059059	0.072499	0.089294	0.1187
PAY_6	0.247234	-0.186558	-0.029488	0.333656	0.402223	0.440377	0.504116	0.668401	1.000000	0.059141	0.071799	0.085883	0.1071

Observations

- The categorical features most highly correlated with 'HeartDisease' are 'Stroke_Yes', 'DiffWalking_Yes', 'AgeCategory_80 or older', 'Diabetic_Yes', 'GenHealth_Fair', 'GenHealth_Poor', 'KidneyDisease_Yes'
- There is not significant collinearity between the features (Rule of thumb: If the correlation > 0.8 then severe multicollinearity may be present)

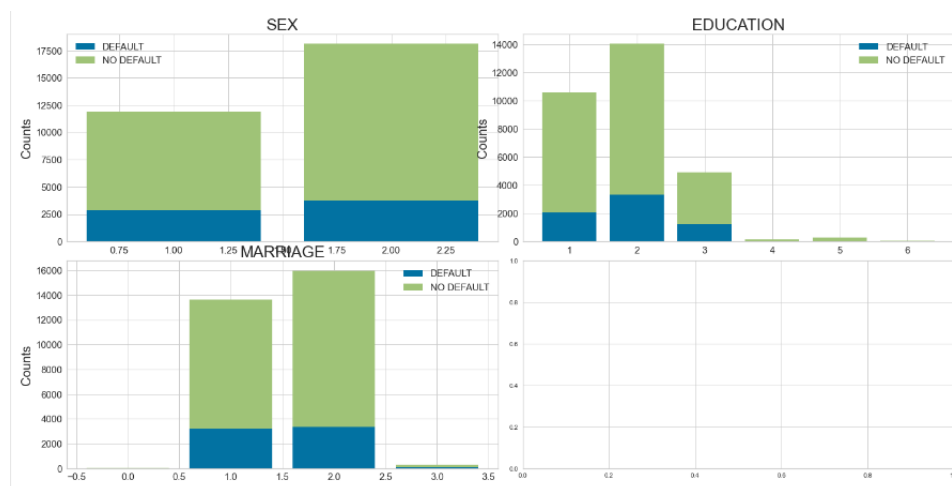
- **Histogram:** Created for only the most correlated features to 'default payment next month' ('SEX', 'EDUCATION', and 'MARRIAGE') to show their distributions and observations were captured:

- Most people in the dataset do not default on their next loan payment
- Most people in the dataset do have a university education (2) followed by graduate education (2) and high school (3)
- Most people in the dataset are female (2)
- Most people in the dataset are single (2) followed by married (1)



- **Stacked Bar Chart:** Created to compare Yes Default and No default by the most correlated features found in the Spearman's Correlation Matrix, and observations were captured:

- Most people in the dataset are female and there are slightly higher loan payment defaults for females because of that
- Most people in the dataset have a university education (2) and there are slightly higher loan payment defaults for females because of that
- There are just as many loan defaults in those who are married as those who are single



- **Train and Test Sets**

- **Train and Test Sets:** To support the training and testing of the models, the dataset was split data into training (80%) and test (20%) sets with 'default payment next month' feature (binary 1=Yes, 0=No) as the target

- **Missing Data:** To ensure that any missing data is captured and then managed, `isna()` was applied to the training and test sets. It was found that there is no missing data in the dataset.
- **5 Best Features:** To find the best 5 features, X2 was used, and additional train & test sets were created using the `chi2` function.
 - The 5 best features from X2 are as follows: `LIMIT_BAL`, `PAY_AMT1`, `PAY_AMT2`, `PAY_AMT3`, `PAY_AMT5`
- **Automatic Class Balancing of the Target ('default payment next month')**
 - **Class Balancing:** "A classification problem in machine learning is where we have given some input (independent variables), and we have to predict a discrete target. It is highly possible that the distribution of discrete values will be very different. Due to this difference in each class, the algorithms tend to get biased towards the majority values present and don't perform well on the minority values. This difference in class frequencies affects the overall predictability of the model." (Singh, 2021). In this dataset, the distribution of the target variable 'default payment next month' is very different (four times as many with 'default payment next month' _No compared to 'default payment next month' _Yes)
 - To remedy this class imbalance, automatic class balancing was used. In automatic class balance logistic model "We have added the `class_weight` parameter to our logistic regression algorithm and the value we have passed is 'balanced'." (Singh, 2021).

- **Automatic class balance** will be implemented into the logistic regression models using `solver='newton-cg', class_weight='balanced'`

Results

- **Models:** A total of 2 Logistic Regression models were created to find the model with the best performance using the loan default dataset. The Logistic Regression models used a combination of all 23 features (ID feature was removed as it is not useful in this analysis), 5 best features from X2, in addition to automatic weighting of the target class. The precision and recall were also calculated, but because the F1 score is a combination of these two, the F1 score is the metric used below.

Model	Accuracy	F1 Score (Yes-Loan Default)	F1 Score (No-Loan Default)	AUC Score
Logistic Regression - Automatic class weighting of target variable with all 23 features	78%	0.522	0.857	0.77
Logistic Regression - Automatic class weighting of target variable with 5 best features using X2	78%	0.526	0.858	0.75

- **Model Findings**
 - **Logistic Regression Models:** 2 logistic regression models were created to find the best model performance with class balancing
 - Both logistic regression models created with class balancing (automatic with all 23 features and automatic with 5 best features from X2) had relatively the same performance (78% accuracy, 0.5 F1 score for minority class (Yes - loan default), 0.8 F1 score for majority class (No – loan default), and 0.75-0.77 AUC score.

- Although the automatic class balancing logistic regression models had relatively the same performance, the fact that the 5 best features from X2 did relatively the same as both all 23 features lends itself to show that these 5 best features (LIMIT_BAL, PAY_AMT1, PAY_AMT2, PAY_AMT3, PAY_AMT5) have a bigger impact on the target variable ('default payment next month') than any of the other variables combined

Conclusion

Multiple Logistic Regression models were used on the dataset to determine which features are most related to or most predictive of loan default, where the 'default payment next month' feature was the target.

- Matrices & Tests used to evaluate features in relation to the target variable, 'default payment next month'

Matrices & Tests	Features
Pearson's & Spearman's Correlation Matrix - Feature correlation to 'default payment next month' target variable	<ul style="list-style-type: none"> • Pearson's Correlation Matrix (numeric features) <ul style="list-style-type: none"> ○ 'PAY_0' ○ 'PAY_2' ○ 'PAY_3' ○ 'PAY_4' ○ 'PAY_5' ○ 'PAY_6' • Spearman's Correlation Matrix (categorical features)

	<ul style="list-style-type: none"> ○ 'SEX_1' (male) ○ 'EDUCATION_2' (university) ○ 'EDUCATION_3' (high school) ○ 'MARRIAGE_1' (married)
Chi-Square (X2) Test - 5 best features correlated to 'default payment next month'	<ul style="list-style-type: none"> ● LIMIT_BAL ● PAY_AMT1 ● PAY_AMT2 ● PAY_AMT3 ● PAY_AMT5

- **Findings:**
 - The Pearson's & Spearman's Correlation Matrix features and the Chi-Squared (X2) features complement each other (with the exceptions of the categorical features), providing more evidence that the X2 5 best features are accurate
 - The 5 best features (chi-squared) of LIMIT_BAL, PAY_AMT1, PAY_AMT2, PAY_AMT3, PAY_AMT5 imply that a loan applicant's Limit Balance and payment history are most influential in predicting if they will default on their next payment
 - With the model performance metrics in mind (see analysis above), the Logistic Regression Models with all 23 Features and only the 5 Best Features performed relatively the same. This proves that the 5 Best Features are in fact as predictive of the target variable as all Features. The Logistic Regression model with 5 best features will be used for analyzing future loan applicants' information and determine if the loan should be approved.

Assumptions

There were a few assumptions of this research initiative, which can be found below:

- Data is accurately provided by banking institutions
- Since a negative credit card balance means that the lender owes the borrower an amount rather than the other way around, changing the negative payment values in PAY_0-6 to 0 to run X2 functions was appropriate. Both a negative and 0 value imply that there was no credit card balance owed.
- Logistic Regression model with 5 Best Features' accuracy of 78% is statistically acceptable
- F1 scores of target variable are acceptable
 - 0.526 for "Yes Default" minority class
 - 0.858 for "No Default" majority class
- AUC value is acceptable at 0.75
 - Excellent for AUC values between 0.9-1
 - Good for AUC values between 0.8-0.9
 - Fair for AUC values between 0.7-0.8
 - Poor for AUC values between 0.6-0.7
 - Failed for AUC values between 0.5-0.6

Limitations

There were some limitations of this research initiative, which can be found below:

- No data of number of dependents, co-applicant information, or credit history score as new features which could make the model's performance better and provide more information to the underwriter about the loan applicant's risk of approving the loan.

Challenges

There were some challenges to this research initiative, which can be found below:

- To have the visualizations x-axes show as words (for example, Male & Female) instead of 1 and 2, these variables were initially replaced with their string equivalents. When doing this however, error messages showed when creating the visualizations. To keep the error messages to a minimum, the categorical variable values were not transformed to a string and remained as a numerical number with comments on what the numbers mean.

Future Uses/Additional Applications

Future applications of this Logistic Regression model with 5 Best Features will be recommended to use in making underwriter decisions on if a loan application should be approved.

If after adding the loan application information into the model and the prediction is "Yes Default" then the application will be rejected. If the prediction is "No Default" then the application will be approved.

Recommendations

As a recommendation, the X2 5 best features are easier to implement in a financial institution environment/setting than all 23 (i.e., it's easier to inform loan underwriters on 5 features that impact or are most related to loan defaults, rather than 23 features). The 5 best features are:

- LIMIT_BAL: Amount of the given credit (NT dollar): it includes both the individual consumer credit and his/her family (supplementary) credit (continuous numerical value in dollars)
- PAY_AMT1: History of past payment (April 2005)
- PAY_AMT2: History of past payment (May 2005)
- PAY_AMT3: History of past payment (May 2005)
- PAY_AMT5: History of past payment (July 2005)

When advising loan underwriter on how to assess how risky customers would be to provide a loan to, these are the recommendations to provide:

- Customers with lower limit balance and duly paid historical monthly loan payments less likely to default on their loans
- These particular customers should be given a better chance of having their loan applications approved as their likelihood of defaulting on that loan is lower than that of customers who have a higher limit balance and pay their historical monthly loan payments late
- Customer incentives for decreasing the risk of loan defaults and increasing their future loan application approval chances could include: free access to credit score websites (e.g.

Equifax, TransUnion) to monitor their loan payments, and free access to money management sites (e.g., Mint.com)

When mortgage brokers are advising clients when applying for loans, these are the 5 recommendations to provide:

- You should keep your limit balance reasonably low
- You should avoid paying your monthly loan payments late
- You should focus on paying your monthly loan payments on time

In future iterations of this loan application initiative, the following are recommended to be implemented:

- Add number of dependents as a new feature
- Add co-applicant information as new features
- Add credit history and/or score as new features
- Create a logistic regression model with 5 Best Features and automatic class balancing (as opposed to automatic which was done in the initial model) based on the dataset instead of automatic class balancing to understand if the model performance improves

Implementation Plan

The logistic model has been trained and then a testing set used to obtain model accuracy, precision and F1 scores. The next step would be to do a real-life test with loan applicant data from financial institutions. The model isn't necessarily ready for deployment, but it is ready to be beta tested with real-life data from loan applications. This will ensure that more potential real-

world scenarios that relate or do not relate to defaulting on the next loan payment are captured in the dataset

To improve the model, financial institution staff (specifically, underwriters) will also be interviewed on loan application data collection and possible methods for improvement. Examples of feedback include different data repositories, new features/measurements from loan applicants to add to the dataset (e.g., number of dependents, number of current financial accounts), possible anecdotal features that underwriters and loan officers have observed in the field.

Ethical Assessment

Ethically, it's important to provide this loan approval prediction information to loan applicants (I.e., what the underwriter is looking at in their application while deciding to approve or reject), but not with the intent to modify the applications with false information just to get approved. This may lead to default loans as the applicant was not able to pay the loan back in the first place, "This is part of the reason we saw a massive housing collapse at the end of last decade. Lenders started giving out loans to people who really weren't qualified, just so they could make a quick buck. But now that we've fully recovered and the economy is as stable as it's been in years, it's time that we study the role of ethics in the lending profession." (Contributor, 2017)

Appendix

- **Accuracy:** “Accuracy represents the number of correctly classified data instances over the total number of data instances” (B, H. N., 2020)
- **Precision:** “Precision talks about how precise/accurate your model is out of those predicted positive, how many of them are actual positive.” (B, H. N., 2020)
- **Recall:** “Recall actually calculates how many of the Actual Positives our model capture through labeling it as Positive (True Positive).” (B, H. N., 2020)
- **F1 Score:** “F1-score is a metric which takes into account both precision and recall.
- **Precision:** positive value
- **Recall:** true positive rate” (B, H. N., 2020)
- **AUC Score:** ““So, what area under the ROC curve describes good discrimination?

Unfortunately, there is no "magic" number, only general guidelines. In general, we use the following rule of thumb:

- 0.5 = This suggests no discrimination, so we might as well flip a coin.
- 0.5-0.7 = We consider this poor discrimination, not much better than a coin toss.
- 0.7-0.8 = Acceptable discrimination
- 0.8-0.9= Excellent discrimination
- >0.9 = Outstanding discrimination" (Dixon, 2020)

References

1. Mariotti, T. (2018, March 21). *6 steps of the Mortgage Loan Process: From pre-approval to closing*. RubyHome.com. Retrieved December 18, 2022, from <https://www.rubyhome.com/blog/mortgage-loan-process/#:~:text=There%20are%20six%20distinct%20phases,loan%20processing%3B%20underwriting%20and%20closing.>

2. Abbott, D. (2014). *Applied predictive analytics: Principles and techniques for the professional data analyst*. Wiley.
3. Ntiamoah, E. (2014). Journal of Accounting, Business and Finance Research. *Loan Default Rate and Its Impact on Profitability in Financial Institutions*.
<https://doi.org/10.20448/2002>
4. Kehoe. (2018, July 4). How to start your first data science project - A practical tutorial for beginners. Hi, I'm Juan L. Kehoe! Retrieved December 4, 2022, from
<https://juan0001.github.io/How-to-start-your-first-data-science-project-a-practical-tutorial-for-beginners/>
5. El Khouli, R. H., Macura, K. J., Barker, P. B., Habba, M. R., Jacobs, M. A., & Bluemke, D. A. (2009, November). Relationship of temporal resolution to diagnostic performance for dynamic contrast enhanced MRI of the breast. Journal of magnetic resonance imaging: JMIR. Retrieved February 20, 2022, from
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2935260/>
6. B, H. N. (2020, June 1). Confusion matrix, accuracy, precision, recall, F1 score. Medium. Retrieved May 14, 2022, from <https://medium.com/analytics-vidhya/confusion-matrix-accuracy-precision-recall-f1-score-ade299cf63cd#:~:text=F1%20Score%20becomes%201%20only,a%20better%20measure%20than%20accuracy.>