

# Machine Learning Analysis of Mortgage Credit Risk

Siva Pillai, Jennifer Woodbury, Megan Polak and Ujwal Pathak  
Akshit Shah and Chaitanya Potnis, Avery Leider and Charles C. Tappert  
Pace University

Pleasantville, NY 10570, USA

Email: {sp57299w, jw82676n, mp32733n, up12137, as87187n, cp39352n, aleider, ctappert}@pace.edu

**Abstract**—In 2008, the US experienced the worst financial crisis succeeding the Great Depression of the 1930s. A recession fueled by an influx of poorly underwritten mortgages, in which a high percentage of “less credit-worthy” borrowers defaulted on their mortgage payments. To date, the market has recovered from the collapse but we must avoid the pitfalls of the previous market meltdown. Greed and over zealous assumptions fueled the 2008 crisis and it is imperative that bank underwriters properly assess risks with the assistance of newer technologies. In this paper we utilize machine learning techniques to predict the approval or denial of a mortgage applicant. The mortgage decision will be determined by a two-tier machine learning model that examines micro and macro risk exposures. We performed comparative analysis using logistic regression, random forest, adaboost optimizer, and deep neural network. Logistic regression provided optimal results and thus the decision model. Our model currently tests at an accuracy level of 85.85% and F1 score of 0.87 using logistic regression. This technology will offer a unique perspective and add value to banking risk models.

**Index Terms**—Machine Learning Model, Mortgage Credit Risk, Logistic Regression, Random Forest Classifier, Deep Neural Network, Classification and Regression Trees, GDP, Unemployment, Home Mortgage Disclosure Act, The Housing and Economic Recovery Act (HERA).

## I. INTRODUCTION

In September 2018, the Board of Governors of the Federal Reserve published the total US mortgage debt outstanding totaling \$15.131 trillion across all holders, an increase of 5% over the previous 12 months. The recent jump sends a strong indication that a recovery has continued in the housing market. In order to prevent another crisis it is imperative that market participants avoid the pitfalls of the US Housing Market melt down in 2008. Mortgage originators and financial institutions practice stricter underwriting guidelines in comparison to the pre-crisis era. Regulators, along with Congress, implemented the Dodd-Frank Wall Street Reform and Consumer Protection Act [1] to assist in preventing a recurrence of the market meltdown from a macro perspective. Regulations to ensure appropriate consumer practices is just the beginning. Additional work is necessary to ensure that borrowers have the ability and commitment to pay their mortgage.

Thanks to the IBM Faculty Award that made this research possible.

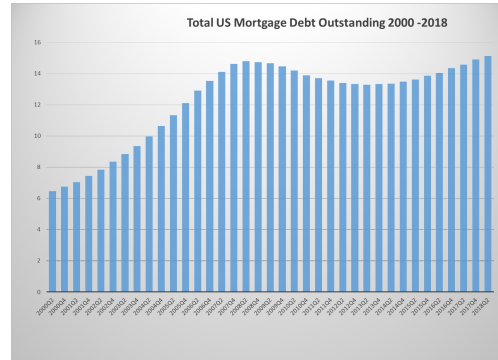


Fig. 1. Source: Board of Governors of the Federal Reserve System (US) Release Date: September 2018

The objective of this project is to build a model that has the ability to assess the credit risk of mortgage related exposures to financial institutions [2]. The model will factor in borrower-level (micro) and market-level stresses (macro) derived while utilizing machine learning technologies. This paper is designed to provide a background of the mortgage market industry, micro and macro level risk exposures, technologies and methodologies used to design and implement the credit risk model. And finally, the results and findings of the machine learning model. [3]

## II. LITERATURE REVIEW

### A. Business Risks and Key Mortgage Elements

1) *Current US Housing Market*: The current state of the US Housing Market reflects a slowdown in popular US regions - Seattle, Silicon Valley and Austin, Texas. Historically, trends in popular US regions set the tone for the market. With rising mortgage rates and prices climbing at a faster rate than income, buyers are getting squeezed and will hit a limit. [4] But market participants continue to view the housing sector as strong due to a healthy labor market and steady economic growth. This indicates price stabilization and not crisis-level conditions. The US Housing market remains strong with interested participants. See Figure 3.

2) *Risk Assessment Models*: Financial institutions rely on proprietary underwriting models to assess their risk exposure to mortgage loans. Underwriters carefully examine

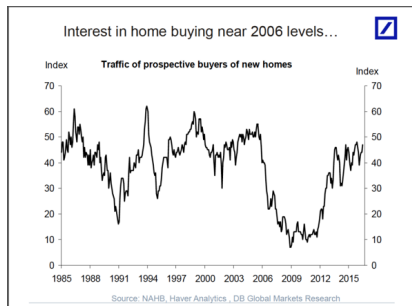


Fig. 2. Source: NAHB, Haver Analytics, Deutsche Bank Global Markets Research

personal information and credit profiles to ensure borrower eligibility. Key elements that play a role in the decision making process include: FICO (Fair Isaac Corporation) score, occupation, total household income, DTI (debt to income) ratio, property location, loan amount, loan to value ratio (LTV), full documentation availability, property type and occupancy status. [5]

These elements are used as cohorts to determine if a borrower qualifies for a mortgage. Researchers conclude that buyers with low LTV ratios and high FICO scores typically qualify for lowest mortgage rates. [4] A borrower with less than pristine credit does not automatically disqualify a borrower. But other factors together may reduce the risk associated with the borrower. For example a borrower with a low LTV (less than 70 percent), high FICO (700+), proof of income/high income, and low DTI (less than 20 percent) is viewed as a less risky profile. [6] The loan amount and LTV ratio combined can be used to determine whether or not private mortgage insurance (PMI) is required.

The first objective of our project will be to utilize key mortgage variables to construct a borrower profile that will not default on their mortgage loan.

### B. Current Market Utilization of Machine Learning

1) "Analysis of Feature Selection Techniques in Credit Risk": Utilizing the best prediction features in credit analysis is crucial in assessing risk. We looked at credit risk assessment to get a better understanding of variables used to assess mortgage credit risk. In Analysis of feature selection techniques in credit risk assessment, R. S. Rama and S. Kumaresan [7] found that the most important features in credit risk assessments are checking account status, credit history, duration in months, saving account balances, purpose, credit score, property type, present employment, occupancy, age, installment plans, personal status, and sex. The feature selection was done using information gain, gain ratio, and chi square correlation. Data used in this research was public data of German credit that consists of 1000 instances in which 700 of them are creditworthy applicants and 300 of bad credit applicants.

| Attribute No | Name of the attribute                                | Value of gain ratio |
|--------------|--|---------------------|
| 1            | Status of checking account                           | 123.7209            |
| 3            | Credit history                                       | 61.6914             |
| 2            | Duration in month                                    | 46.8311             |
| 6            | Savings account                                      | 36.0989             |
| 4            | Purpose  | 33.3564             |
| 5            | Credit amount  | 26.9528             |
| 12           | Property   | 23.7196             |
| 7            | Present employment since                             | 18.3683             |
| 15           | Housing  | 18.1998             |
| 13           | Age in year  | 16.3681             |
| 14           | Other installment plans                              | 12.8392             |
| 9            | Personal status and sex                              | 9.6052              |
| 20           | Foreign worker                                       | 6.737               |
| 10           | Others debtors                                       | 6.6454              |
| 17           | Job  | 1.8852              |
| 19           | Telephone  | 1.3298              |
| 18           | Number of people being liable to provide maintenance | 0                   |
| 8            | Installment rate in percentage                       | 0                   |
| 11           | Present residence since                              | 0                   |
| 16           | Number of existing credits at this bank              | 0                   |

Fig. 3. Source: R. S. Ramya and S. Kumaresan, "Analysis of feature selection techniques in credit risk assessment," 2015 International Conference on Advanced Computing and Communication Systems, Coimbatore, 2015, pp. 1-6.

2) "A Machine Learning Approach for Predicting Bank Credit Worthiness": Analysis of UCI machine learning dataset revealed that there is a relationship between the customer's age and their account balance. Customers who are between 20 and 60 years of age and have small bank account balance are most prone to become defaulters. The dataset contained 23 variables from which researchers picked 5 most important features. Using the 5 selected features they performed multiple classifications. Each of these algorithms achieved an accuracy rate between 76% to over 80%. [8] See Fig 6 and Fig 7.

### C. Market Risk Factors

1) GDP: Gross domestic product (GDP) measured quarterly and annually, provides insight into the growth rate of a nation's economy. GDP is measured as nominal GDP (inflation included) and real GDP (excludes inflation). [9]. Two major factors that affect GDP include: inflation and recession. The Federal Reserve uses 2 policies to maintain GDP and the Economy. **Expansionary Monetary Policy** to ward off recession and **Contractionary Monetary Policy** to prevent inflation. Both policies have a major effect on the disposable income of American households. Under expansionary policy, interest

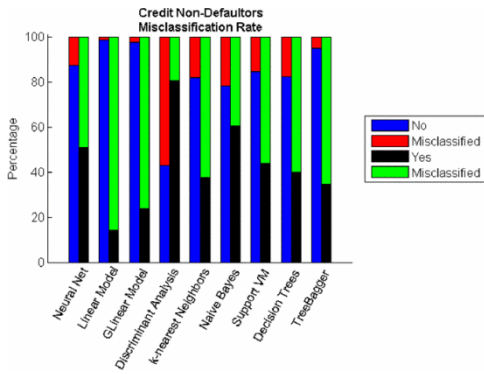


Fig. 4. Source: R. E. Turkson, E. Y. Baagyere and G. E. Weny, "A machine learning approach for predicting bank credit worthiness," 2016 Third International Conference on Artificial Intelligence and Pattern Recognition (AIPR), Lodz, 2016, pp. 1-7.

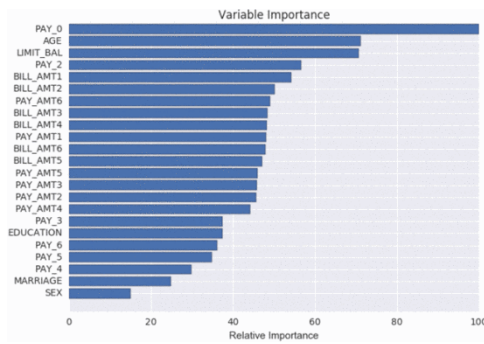


Fig. 5. Source: R. E. Turkson, E. Y. Baagyere and G. E. Weny, "A machine learning approach for predicting bank credit worthiness," 2016 Third International Conference on Artificial Intelligence and Pattern Recognition (AIPR), Lodz, 2016, pp. 1-7.

rates are lowered making it cheaper to borrow and reducing the incentive to save. While contractionary policy aims to decrease the money supply by reducing price levels, and increase private consumption.

The US Federal Reserve utilizes certain tools to maintain GDP and the US Economy: open market operations, discount rate, and reserve requirements.

- Open Market Operations: Central Banks buy and sell securities in the open market.
- Discount Rate: The rate that Central Banks charge its members to borrow at its discount window.
- Reserve Requirement: The required money that banks are mandated to hold overnight.

2) *Unemployment*: Definition: Unemployment occurs when a person who is actively searching for employment is unable to find work. Unemployment is often used as a measure of the health of the economy. The most frequent measure of unemployment is the unemployment rate, which is the number of unemployed people divided by the number of people in the labor force. [10]

Effect of Unemployment to the Economy:



Fig. 6. Source: United States Gross Domestic Product 1960 - 2018

- Unemployment causes an individual to be deprived of resources that trickle down to the benefit of society.
- The economy produces 70 percent which contributes to direct consumption. If people begin losing their jobs, the whole cycle will be hampered. As a result GDP is reduced and the country drifts away from making efficient allocation of the resources.
- Unemployment triggers inflationary conditions causing a rise in the general price of goods and services, while the purchasing power of currency decreases.

### III. METHODOLOGY

#### A. Exploratory Data

The Housing and Economic Act of 2008 (HERA) [11] stipulates that certain mortgage information must be made publicly available and stored in a public use database. FHLB adheres to this rule by storing census-level data relating to mortgages purchased by the organization.

For the purpose of this project, we extracted data from FHLB's Public Use Database (PUDB) [12] for mortgage loans acquired from 2010 to 2017 to perform exploratory analysis for creating a base line credit profile. Key fields extracted from the database include: Year, FIPStateCode, FIPSCountyCode, Income, IncomeRatio, UPB, LTV, MortDate, Purpose, Product, Term, AmortTerm, Front, Back, BoCredScor. Additional fields were derived: State, County (State and County were mapped from the US Census Bureau) [13] and PMT (derived from Rate, AmortTerm, and Amount).

#### B. Data Input for ML Model

In order to improve the accuracy results of the machine learning model, additional data was needed that displayed both approved and declined mortgage applicants. As a primary source of data, we extracted 2009 - 2017 annual data reported by financial institutions required by the Home Mortgage Disclosure Act (HMDA) [14]. In 1975, the United States Congress enacted a regulation that required financial institutions to track and ensure fair lending practices

| Column         | Definition   |
|----------------|--|
| Year           | Year Loan was reported   |
| FIPSStateCode  | FIPS State Code  |
| FIPSCountyCode | FIPS County Code   |
| Income         | Total Borrower(s) Annual Income in Whole Dollars   |
| Incrat         | Borrower Income Ratio  |
| UPB            | Acquisition Unpaid Principal Balance in Whole Dollars  |
| LTV            | Loan to Value Ratio at Origination   |
| MortDate       | Year of Mortgage Note  |
| Purpose        | Loan Purpose 1=purchase; 2=refinancing; 3=second mortgage; 4=new construction; 5=rehabilitation                |
| Product        | Product Type 01=Fixed Rate; 02=ARM; 03=Balloon; 04=GPM/GEM; 05=Reverse Annuity Mortgage; 06=other              |
| Term           | Term of Mortgage at Origination in months  |
| AmorTerm       | Amortization Term in months  |
| NumBor         | Number of Borrowers  |
| Occup          | Occupancy Code 1=Principal residence/owner-occupied; 2=second home; 3=investment property (rental)             |
| Rate           | Interest Rate  |
| Amount         | Loan Amount in Whole Dollars   |
| Front          | Front-end Ratio  |
| Back           | Back-end Ratio   |
| BoCreditScor   | Credit Scores are separated into ranges: 1 = < 660, 3 = 660 < 700, 4 = 700 < 760, 5 = 760 or greater 9=missing |
| PMT            | Monthly mortgage payment (derived)   |
| State          | Derived from FIPS State Code   |
| County         | Derived from FIPS County Code  |

Fig. 7. Metadata defining Exploratory Data. Source: <https://www.ffiec.gov/hmda/nationalarchives.htm>

throughout the United States.

The focal fields extracted from this data set include: Action Taken Type, Year, Loan Type, Loan Purpose, Property Type, Occupancy, Amount, State Code, County Code, Income, Denial Reason, Purchaser Type.

|                    |  |
|--------------------|--|
| Year               |  |
| State              | Two-digit FIPS state identifier                                    |
| County             | Three-digit FIPS county identifier                                 |
| Loan Type          | 1 -- Conventional (any loan other than FHA, VA, FSA, or RHS loans) |
| Property Type      | 1 -- One to four-family (other than manufactured housing)          |
| Loan Purpose       | 1 -- Home purchase   |
| Owner-Occupancy    |  |
| Loan Amount        | in thousands of dollars  |
| Action Taken       |  |
| Approved           | 1 -- Loan originated   |
| Approved           | 2 -- Application approved but not accepted                         |
| Declined           | 3 -- Application denied by financial institution                   |
| Approved           | 6 -- Loan purchased by the institution                             |
| Declined           | 7 -- Preapproval request denied by financial institution           |
| Income             | Gross Annual Income  |
| Reasons for Denial |  |
|                    | 1 -- Debt-to-income ratio  |
|                    | 2 -- Employment history  |
|                    | 3 -- Credit history  |
|                    | 4 -- Collateral  |
|                    | 5 -- Insufficient cash (downpayment, closing costs)                |
|                    | 6 -- Unverifiable information                                      |
|                    | 7 -- Credit application incomplete                                 |
|                    | 8 -- Mortgage insurance denied                                     |
|                    | 9 -- Other   |

Fig. 8. Metadata defining HMDA Data. Source: <https://www.fhfa.gov/DataTools/Downloads/Pages/FHLBank-Public-Use-Database-Previous-Years.aspx>

### C. Machine Learning Techniques

The initial phase of the project was dedicated to deriving a baseline profile, stratified by state and application year. The extracted FHLB data is distributed in multiple panda dataframes to cleanse NaN values, and remove unwanted data. The next steps include data exploration, in which identified variables are plotted on a graph (Matplotlib, Seaborn) to determine correlation and highlight key variables that are most impactful to the outcome. Final exploratory steps include normalizing and preprocessing the remaining data. The outcome serves as the basis for baseline assumptions.

The second data set extracted from HMDA serves as primary input for the Classification Model. The categorical data details approved and declined mortgage transactions from 2009 -2017. A logistic regression is then performed on

the data to determine the binary outcome of the model.

Before preprocessing, both datasets are stored in buckets on Google Cloud Storage Browser.

The accuracy of the model will be calculated using cross entropy and according to the loss obtained, the weights and bias will be adjusted to obtain a higher accuracy. Success will be defined as: The loan was approved after risk assessment. Failure will be defined as:

- The loan was approved by the bank, but was declined by the applicant due to a better alternative.
- The application was rejected by the bank because the applicant did not meet the criteria for approval.

## IV. RESULTS

### A. Data Processing - Preprocessor 1.5

As an initial step the Raw Data is cleaned and normalized in a Jupyter Notebook termed "Preprocessor 1.5". The preprocessor utilizes python script to stratify and analyze the mortgage data extracted from FHLB public database.

The cleaning process extracts all loans that fall outside the set criteria. For example loan balances less than \$10000, Income Ratios less than 0.01 and equal to 1. There were approximately 1500 loans deleted from the raw data set. After the loan extractions the data is ready to be analyzed.

Key fields were selected and statistical analysis performed to determine the mean, median, mode, distribution and standard deviation. These same fields were grouped by distribution and the mean or mode is used to determine the base case model. See Figure 4.

### B. Preprocessor for Machine Learning Model

Data was extracted from Home Mortgage Disclosure Act (HMDA) website. We performed feature engineering on the raw data in a separate Jupyter Notebook termed "Preprocessor\_ML.ipynb". The data was processed as follows:

- Loan\_Type filtered for 30 year Conventional Loans.
- Property\_Type filtered for one-to-four family dwellings
- Statecode mapped to State's abbreviation.
- Rate mapped to the state's 30 year Fixed Rate Mortgage (FRM) for the year of mortgage application else used interest rates provided.
- PMT - monthly mortgage payment derived with pmt pandas function; parameters include: loan amount, rate, and loan amortization term.
- GDP - GDP rate mapped for each state by year.
- IncomeRatio - Derived from dividing monthly payment by monthly gross income.
- Credit Score - assumed failing Credit Score if reason for denial included credit history, else credit score assumed to be passing.
- LTV - LTV is provided a pass or fail based on loan denial reason. If "mortgage insurance denied" was reason for the

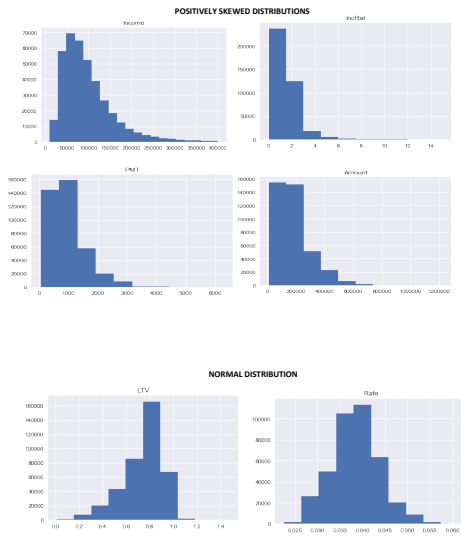


Fig. 9. Field Distributions Positively Skewed and Normal Distributions

loan’s denial, the LTV was assigned 100% else LTV was assigned 75%.

All fields were extracted and saved as "ML\_Processed\_Data.csv". This csv file serves as the primary input for the Classification Model - Machine Learning Model.

### C. CART MODEL- Baseline Model Determination

The model is built on a balanced set of training and testing data. Classification and Regression Tree (CART) algorithm [15] is used to predict the approval or denial outcome. A decision tree will be formed, where each root node represents a variable input and split point on the variable. The leaf nodes are the output variables that will be tabulated to form a final prediction score. Fig 3 illustrates the project’s decision tree algorithm used to calculate the applicant’s prediction score. If the prediction score is greater than 30, it is highly likely that the applicant will be approved, otherwise it will be denied. Variable inputs and splits are:

- Is your income greater than \$10,000?
- Is your Borrower Credit Score greater than 700?
- Is your Income Ratio less than 20%?
- Is the loan amount greater than \$100,000?
- Is the Loan-To-Value (LTV) less than 75%?

### D. CLASSIFICATION MODEL MACHINE LEARNING:

ML\_Processed\_Data.csv is initially read and stored as a dataframe. CSV file generated from Preprocessor\_Machine Learning Model.

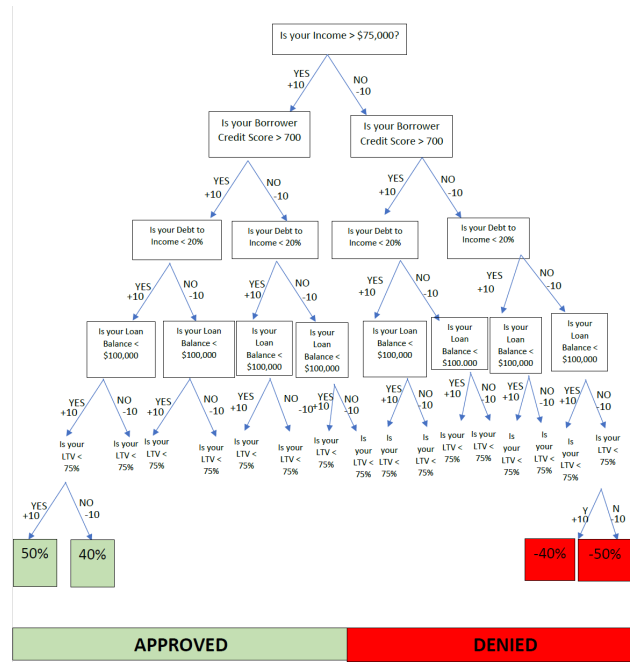


Fig. 10. Decision Tree Used to Determine Prediction Score

1) *Label or Target Creation:* Action Type drives the approved or declined labeling process as follows:

- Approved : Action Type 1, 2, 6 – Label = 1
- Declined: Action Type 3, 7 – Label = 0

2) *Filtering and Selection Fields:* The following interested fields were selected to train the model:

- interested = AgencyCode, LoanType, PropertyType, LoanPurpose, Occupancy, Amount, ActionType, StateCode, CountyCode, Income, PurchaserType, ApplicationDateIndicator, PropertyLocation, USPSCode, GDP, RealState-Growth%, Rate, PMT, IncRat, Unemployment, AmorTerm, BoCreditScor, LTV

3) *Extracting Categorical and Continuous Features:* Interested fields were classified as categorical or continuous categories as below:

- categorical = Agency\_Code, Loan\_Type, Property\_Type, Loan\_Purpose, Occupancy, USPS\_Code, County\_Code, BoCreditScor, LTV
- continuous = Amount, Income, GDP, RealStateGrowth%, Rate, PMT, IncRat, Unemployment

4) *One Hot Encoding:* Applied One Hot encoding on the categorical features

5) *Normalize Continuous Features:* Normalized the continuous features to return an outcome in the range of 0 to 1.

6) *Balance Dataset*: For optimal results the dataset needs to be balance, preferably 50% - 50% split.

7) *Train Test Split*: Total dataset contains 14,000 records and split: Train Data = 2014, 2015, 2016 (288,000 Sample)  
Test Data = 2017 (80,855 Sample)

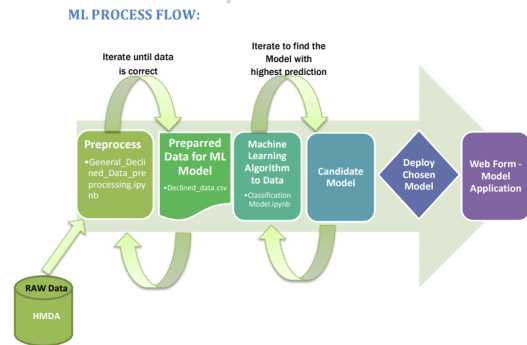


Fig. 11. Machine Learning Process Flow

8) *TRAIN MACHINE LEARNING MODEL*: After performing data pre-processing, a classification algorithm is used to classify the elements into two groups (Approved = 1 and Declined = 0).

Logistic Regression (LR) is used as an efficient machine learning model. LR model is trained on Train data set and the trained model is tested with Test data set and an accuracy of 85.85% is achieved. Figure 6 represents the Confusion Matrix which displays the performance of the algorithm.

|                  |          | Predicted     |          |       |
|------------------|----------|---------------|----------|-------|
|                  |          | Declined      | Approved |       |
| Actual           | Declined | 10370         | 3619     | 13989 |
|                  | Approved | 7825          | 59041    | 66866 |
|                  |          | 18195         | 62660    |       |
| <b>Accuracy:</b> |          | <b>85.85%</b> |          |       |

Fig. 12. Confusion Matrix Logistic Regression

9) *CANDIDATE MODELS*: In total, 9 supervised classification models were tested on the data set. The top four models with the high performance results include: Logistic Regression, Random Forest Classifier, AdaBoost Classifier and Deep Neural Network.

All of the chosen parameters were selected after tuning and testing, performing hyper parameter tuning, cross-validation, etc and the featured parameters are the best for our data selection. Each model parameters are presented below:

Model 1 Logistic Regression Hyper Parameters:

- C =1000,
- max\_iter =100,

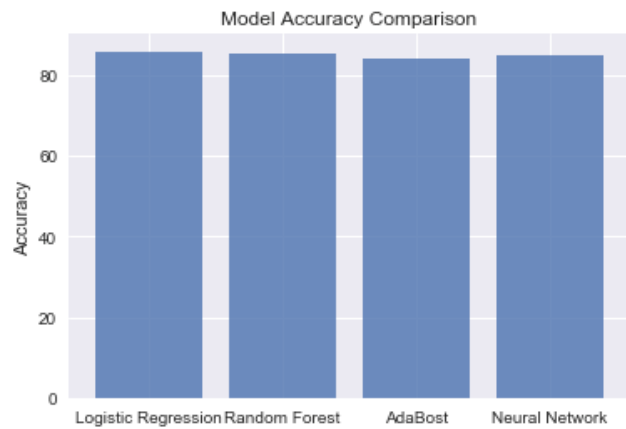


Fig. 13. Model Accuracy Comparison

- solver =liblinear,
- penalty = l2

Model 2 Random Forest Classifier Hyper Parameters:

- n\_estimators =100,
- criterion = 'entropy'

Model 3 AdaBoost Classifier Hyper Parameters:

- algorithm='SAMME.R',
- base\_estimator=None,
- learning\_rate=1.0,
- n\_estimators=50

Model 4 Deep Neural Network

- Layer 1 : 32 units Dense Layer with activation relu
- Layer 2 : 64 units Dense Layer with activation relu
- Layer 3 : 1 units Dense Layer with activation sigmoid
- Loss Function : binary\_crossentropy
- Optimizer : ADAM
- Epochs: 5
- Batch Size : 64
- Validation split : 0.20

10) *CANDIDATE MODEL RESULTS - MODEL SELECTION*: The accuracy levels of each model range between 84% - 85% and all F1-Scores are above 0.85. Accuracy results appear to be similar across models but training time will ultimately be the differentiating factor.

## V. PROJECT REQUIREMENTS

### A. Programming Language

Python will be the programming language utilized throughout the project. All project code is stored in Jupyter Notebooks.

|           | Logistic Regression | Random Forest Classifier | AdaBoost Classifier | Deep Neural Network |
|-----------|---------------------|--------------------------|---------------------|---------------------|
| Accuracy  | 85.84%              | 85.35%                   | 84.15%              | 84.82%              |
| Precision | 0.88                | 0.89                     | 0.87                | 0.87                |
| Recall    | 0.86                | 0.85                     | 0.84                | 0.84                |
| F1-Score  | 0.87                | 0.86                     | 0.85                | 0.85                |

Fig. 14. Candidate Model Results

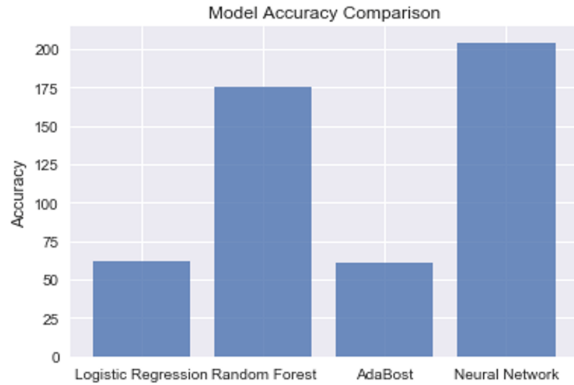


Fig. 15. Training Time Model Chart

#### B. Machine Learning Framework and Libraries:

- Scikit Learn Library
- Keras (with Tensorflow as Backend)

#### C. Google Cloud Storage:

The data in .csv format is uploaded to a google cloud storage bucket called "capstone-project". To import the data into pandas dataframe, the code is as follows:

```
import pandas as pd
import gcsfs
```

```
fs = gcsfs.GCSFileSystem(project='Capstone2')
with fs.open('capstone692/2017_PUDB_Export.csv') as f:
df = pd.read_csv(f)
```

#### D. Data Processing Framework and Libraries:

- Pandas
- Numpy

#### E. Visualization:

- Matplotlib
- Seaborn
- Plotly
- Tableau

|                         | Logistic Regression | Random Forest Classifier | AdaBoost Classifier | Deep Neural Network |
|-------------------------|---------------------|--------------------------|---------------------|---------------------|
| Training Time (Seconds) | 64.01               | 174.85                   | 61.19               | 204.06              |

Fig. 16. Training Time Model Results

## VI. CONCLUSION

After training and testing the ML model to obtain high accuracy we concluded that a single baseline across all states would not be optimal. Instead we implemented a single variable selection process to create base line models for each state and vintage year of mortgage origination.

In addition we discovered HMDA data was the best source of public data that provided approved and declined decisions. After preprocessing and feature engineering, a higher accuracy level was achieved. The model has been trained on 2014 -2016 mortgage data. For test purposes we utilized 2017 data, not used to train the machine model. The model results returned 85.45% accuracy for 2017 mortgage data.

**Finally the Model of Choice is the Logistic Regression, 64.01 seconds to train the model with the highest accuracy of 85.85% and highest F1 score of 0.87.**

All research pertaining to this paper can be found on github at: <https://github.com/JWoodbury125/machinelearningmtgecreditrisk>

## REFERENCES

- [1] SEC Emblem, *Implementing the Dodd-Frank Wall Street Reform and Consumer Protection Act*, 2018. <https://www.sec.gov/spotlight/dodd-frank.shtml>.
- [2] K. MERTENS and M. O. RAVN, "The macroeconomic effects of government asset purchases: Evidence from postwar us housing credit policy andrew i. fieldhouse," *The Quarterly Journal of Economics*, vol. 1, p. 58, 2018.
- [3] W. J. Hippler, S. Hossain, and M. K. Hassan, "Financial crisis spillover from wall street to main street: further evidence," *Empirical Economics*, pp. 1-46, 2018.
- [4] T. Critchfield, J. Dey, N. Mota, S. Patrabansh, *et al.*, "Mortgage experiences of rural borrowers in the united states: Insights from the national survey of mortgage originations," tech. rep., Federal Housing Finance Agency, 2018.
- [5] L. Guiso, A. Pozzi, A. Tsoy, L. Gambacorta, and P. E. Mistrulli, "The cost of distorted financial advice: Evidence from the mortgage market," 2018.
- [6] G. Calcagnini, R. Cole, G. Giombini, and G. Grandicelli, "Hierarchy of bank loan approval and loan performance," *Economia Politica*, pp. 1-20, 2018.
- [7] R. Ramya and S. Kumaresan, "Analysis of feature selection techniques in credit risk assessment," in *Advanced Computing and Communication Systems, 2015 International Conference on*, pp. 1-6, IEEE, 2015.
- [8] R. E. Turkson, E. Y. Baagyere, and G. E. Wanya, "A machine learning approach for predicting bank credit worthiness," in *Artificial Intelligence and Pattern Recognition (AIPR), International Conference on*, pp. 1-7, IEEE, 2016.
- [9] J. N. Inekwe, Y. Jin, and M. R. Valenzuela, "The effects of financial distress: Evidence from us gdp growth," *Economic Modelling*, vol. 72, pp. 8-21, 2018.

- [10] R. S. Kaplan *et al.*, “Discussion of economic conditions and key challenges facing the us economy,” tech. rep., Federal Reserve Bank of Dallas, 2018.
- [11] B. Arthur, “Housing and economic recovery act of 2008,” *Harv. J. on Legis.*, vol. 46, p. 585, 2009.
- [12] Federal Home Loan Bank, *Fhfa.gov. (2018). Federal Home Loan Bank Member Data —Federal Housing Finance Agency*, 2018. <https://www.fhfa.gov/DataTools/Downloads/Pages/Federal-Home-Loan-Bank-Member-Data.aspx>.
- [13] United States Census Bureau, *FIPS Codes for the States and the District of Columbia*, 2018. [https://www.census.gov/geo/reference/ansi\\_s\\_tatetables.html](https://www.census.gov/geo/reference/ansi_s_tatetables.html).
- [14] Federal Financial Institutions Examination of Council, *Home Mortgage Disclosure Act*, 2018. <https://www.ffiec.gov/hmda/default.htm>.
- [15] D. G. T. DENISON, B. K. MALLICK, and A. F. M. SMITH, “A bayesian cart algorithm,” *Biometrika*, vol. 85, no. 2, pp. 363–377, 1998.