

Kyle Nero, Chung En (Shawn) Pan, Nathan Van Lingen, Koosha Jadbabaei

Industry Mentor: Brian Duke (Petal)

Faculty Mentor: Berk Ustun

Website Link

Abstract

In this study, we classify banking transactions for the purposes of credit scoring and credit applications. These classified transactions could help us assess credit worthiness in the absence of a credit history. For people like immigrants, students, or anyone looking to build credit and is without a good history, this is a convenient metric that is available for nearly any credit applicants. We evaluate various methods for data cleaning, extraction, and modeling, and build an ensemble model to further improve performance and reduce overfitting.

Introduction: Traditional Credit Card Application Process

Applicants fill out a loan application form for a credit card or other personal loan Applicants submit the form along with required documents* to the lenders for review Lenders evaluate applicants' credit worthiness (credit scores) using credit scoring models

High credit scores gained pre-approval, but banks checks documents for final approval

Application Submission Review Approval

Required documents may contain proof of identity such as SSN, as well as income verification

Review Approval

Low credit score leads to rejection.

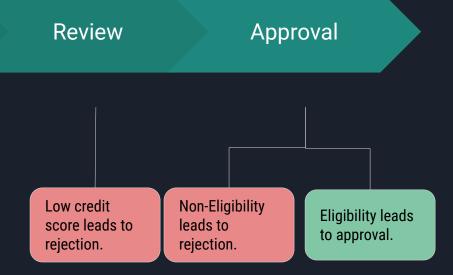
Review Approval

Eligibility leads to rejection.

Introduction: Traditional Credit Card Application Process

Applicants fill out a loan application form for a credit card or other personal loan Applicants submit the form along with required documents* to the lenders for review Lenders evaluate applicants' credit worthiness (credit scores) using credit scoring models High credit scores gained pre-approval, but banks checks documents for final approval

Applicants who are foreigners, college students, and individuals who did not have credit score often receive rejections because they did not have a strong basis of credit score!



Introduction: How is Petal Changing the Game

Applicants fill out a loan application form for a credit card or other personal loan Applicants submit the form along with required documents* to the lenders for review Lenders evaluate applicants' credit worthiness (credit scores) using credit scoring models

High credit scores gained pre-approval, but banks checks documents for final approval



CHEQUING ACCOUNT STATEMENT
Page: 1 of 1

OHN JONES	Statement period	Account N	
643 DUNDAS ST W APT 27	2003-10-09 to 2003-11-08	00005-	
ORONTO ON M6K 1V2		123-456-	

Date	Description	Ref.	Withdrawals	Deposits	Balance
2003-10-08	Previous balance				0.55
2003-10-14	Payroll Deposit - HOTEL			694.81	695.36
2003-10-14	Web Bill Payment - MASTERCARD	9685	200.00		495.36
2003-10-16	ATM Withdrawal - INTERAC	3990	21.25		474.11
2003-10-16	Fees - Interac		1.50		472.61
2003-10-20	Interac Purchase - ELECTRONICS	1975	2.99		469.62
2003-10-21	Web Bill Payment - AMEX	3314	300.00		169.62
2003-10-22	ATM Withdrawal - FIRST BANK	0064	100.00		69.62
2003-10-23	Interac Purchase - SUPERMARKET	1559	29.08		40.54
2003-10-24	Interac Refund - ELECTRONICS	1975		2.99	43.53
2003-10-27	Telephone Bill Payment - VISA	2475	6.77		36.76
2003-10-28	Payroll Deposit - HOTEL			694.81	731.57
2003-10-30	Web Funds Transfer - From SAVINGS	2620		50.00	781.57
2003-11-03	Pre-Auth. Payment - INSURANCE		33.55		748.02
2003-11-03	Cheque No 409		100.00		648.02
2003-11-06	Mortgage Payment		710.49		-62.47
2003-11-07	Fees - Overdraft		5.00		-67.47
2003-11-08	Fees - Monthly		5.00		-72.47
	*** Totals ***		1,515.63	1,442.61	

Review

Approval

To avoid the problem of the "credit invisible" not easily getting approval for credit, Petal allows applicants to submit optional checking account statements as a metric of evaluation for the banks. They use this data, including transaction date, amount, and memo to create new features for their credit models. To process this data, these transactions must be grouped into different categories (like food, automotive, atm, health, etc.), which is what we aim to do in this project

Eligibility leads to approval.

Creating the Perfect Dataset

Text Cleaning

- Removing Stop Words (like a, an, the, and, is)
- Removing all punctuation and spaces
- Made all words lowercase
- Ex: "Wal-Mart" → "walmart"

memo
POS CASINO BAR @ SPOTL - MEMO=PURCHASE 03/02 C
BEST BUY GRAND REGENCY BRANDON F
CORNER STORE ARLINGTON TX 10/17 Purchase \$5.3
SPEEDWAY IN BEDFORD IN 07/10 DEBIT_CARD
PAYMENT FOR AMZ STORECARD WEB ID: ACH_DEBIT

pos casino bar spotl memopurchase 0302 coache...

best buy grand regency brandon f

corner store arlington tx 1017 purchase \$536 ...

speedway bedford 0710 debitcard

payment amz storecard web id achdebit

Feature Engineering

- Applying TF-IDF to text
- Standardized \$ amount
- Created new columns like
 - Is amount a whole number?
 - Was purchase made on a holiday? A weekend?
- One Hot Encoded
 - Year
 - Month
 - Day

	transaction_date	amount	
8	2018-03-06	15.25	
39	2018-09-29	16.04	
45	2018-10-17	10.36	
52	2017-07-11	4.63	
55	2018-09-10	223.00	

	amount (standardized)	is_whole_number	is_holiday	is_weekend
8	-0.012686	0	0	0
39	-0.012246	0	0	1
45	-0.015409	0	0	0
52	-0.018601	0	0	0
55	0.103030	1	0	0

EDA: Cleaned Word Clouds



Food and Beverages: 162,009 (32.66%)



General Merchandise: 138,961 (28.01%)



Automotive: 63,617 (12.82%)

```
wal walmart was super center; a super center;
```

Groceries: 63,541 (12.81%)



Entertainment: 31,492 (6.35%)



Travel: 25,491 (5.14%)



Healthcare / Medical: 8,392 (1.69%)

```
petsmart memopurchase
animal authorized petsmart supplies properties properti
```

Pets / Pet Care: 2,588 (0.52%)

Model Pipeline



Text Cleaning/Feature Engineering

Applied Logistic Regression Classifier



BERT:

Poor accuracy due to lack of syntax in memos

- Applied XGBoost on newly engineered non-text features
- ~ 34% Accurate

Baseline Non-Text Model:

Tuned Text-Only Model:

Baseline Text-Only Model:

on TF-IDF features ~ 85% Accurate

- Tested both word and char grams and n-gram sizes with Logistic Regression
- ~88% Accurate

Hardcode Text-Only Model:

Found ~1000 keywords that correlated with a given category 90%+ of the time



Combined Text/Non-Text Models:

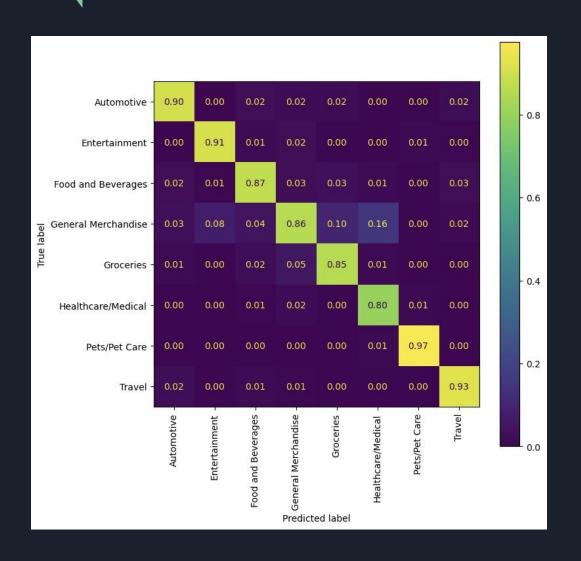
- Creating ensemble model to combine tuned logistic regression models with tf-idf and other engineered features
- ~89% Accuracy

Combined Text-Only Models:

- Created ensemble model to combine tuned logistic regression model and hardcoded model
- ~86% Accurate (lower accuracy)
- No longer using hardcoded model

Final Model Performance

- Best Model
- Tuned Text Model (TF-IDF/Logistic Regression) and Non-Text Model (XGBoost) Ensemble → ~89% Accuracy



Category	Accuracy	
Automotive	90%	
Entertainment	91%	
Food and Beverages	87%	
General Merchandise	86%	
Groceries	85%	
Healthcare/Medical	80%	
Pets/Pet care	97%	
Travel	93%	

Conclusion

Despite credit score being a significant metric for lenders during the loan application approval process, this single metric fails to consider how underrepresented demographics may be deemed as unworthy of credit when in fact they should be worthy. This is a lose-lose for both lenders and applicants. On one hand, lenders lose potential customers ("the credit invisible") who are excluded due to this traditional process. At the same time, these applicants are not allowed to receive credit, which makes things like buying a house seemingly impossible. We aim to explore the limitations of traditional credit score models and propose an alternative method for determining creditworthiness that is more inclusive and equitable. We approach this problem by utilizing supplemental features like a user's categorized transaction history to the traditional credit scoring model. With a user's bank statements, we could extract information like the transaction date, amount, and memo to flag each transaction into a category (with ~89% accuracy). The next step would be to utilize these categories in creating new features to optimize common credit scoring models to make them stronger and more fair. This approach solves both of the problems we described above as it will aid applicants with low/no credits as well as profit-hungry institutions looking to acquire more customers in financial industry.

Appendix

Bank statement Examplehttps://en.wikipedia.org/wiki/Bank statement

