

Credit Default Risk Analysis

Lina Gao

Problem Statement

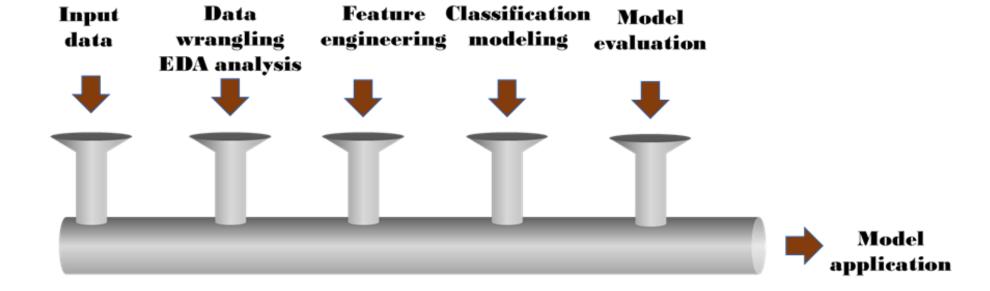
• How do we use modern machine learning algorithms to identify potential credit defaulter based on their historical transaction data and socioeconomic status?

Source of Data

- ➤ It includes 30000 observations with 25 features collected in Taiwan from April 2005 to Sep 2005.
- Features: Age, Sex, Marriage, Education, credit limit, payment, bill and payment status in each of the six months.



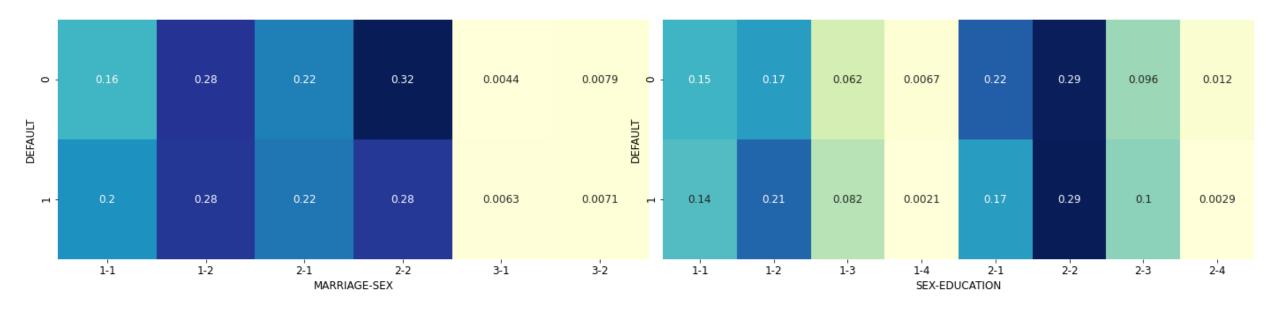
Analytical Workflow



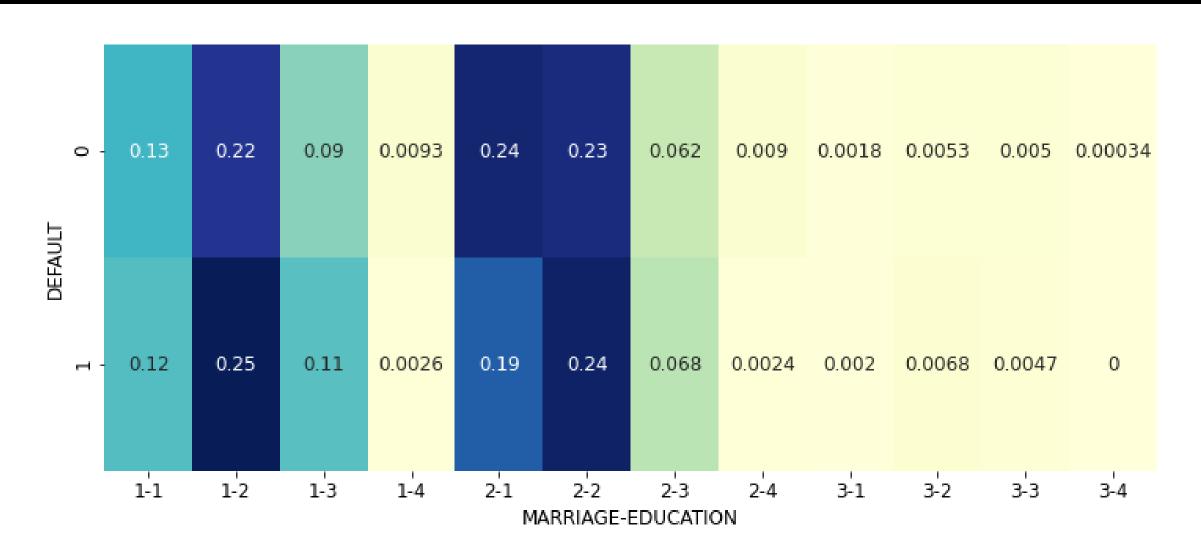
Data structure

Variables	Description
ID	ID of each client
LIMIT_BAL	Amount of given credit in NT dollars (includes individual and family/supplementary credit)
MARRIAGE	1=married, 2=single, 3=other
SEX	Gender (1=male, 2=female)
EDUCATION	(1=graduate school, 2=university, 3=high school, 4=others, 5=unknown, 6=unknown)
DEFAULT	1=default, 0=non-default. * N(non-default):N(default)=3.52:1
PAY_0, PAY_2,	Repayment status in September, August, July, June, May, April (-2=no consumption, -1=pay duly,
, PAY_6	0=the use of revolving credit, 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,	Amount of bill statement in September, August, July, June, May, April, 2005 (NT dollar)
BILL_AMT2,	
••••	
BILL_AMT6	
PAY_AMT1,	Amount of previous payment in September, August, July, June, May, April 2005 (NT dollar)
PAY_AMT2,	
••••	
PAY_AMT6	

Data Exploration: client segmentation



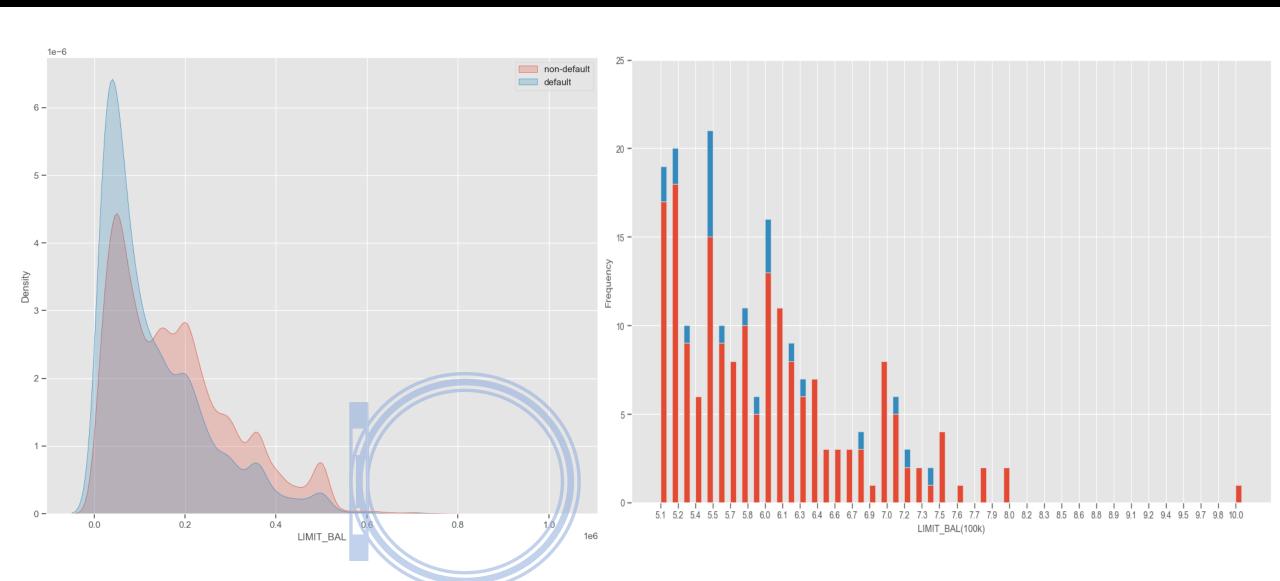
Data Exploration



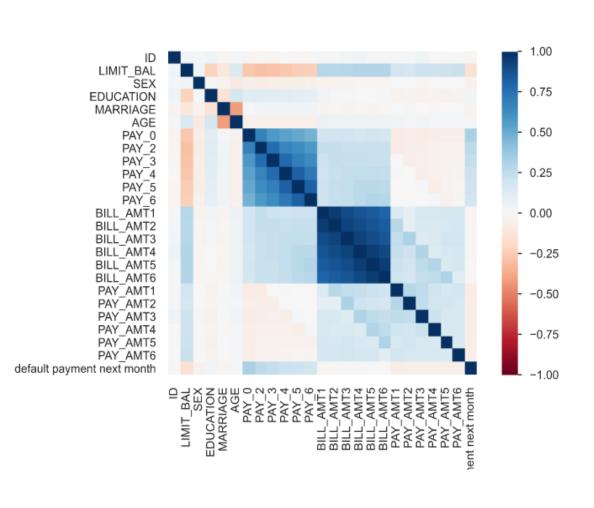
Statistics test to explore the relationship between the features

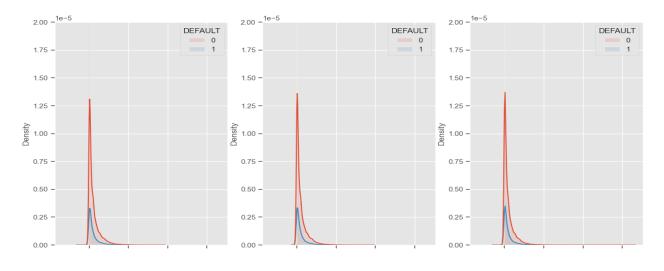
features	Statistic tests	Results
DEFAULT vs SEX	Chi-square test	P<0.05
DEFAULT vs MARRIAGE		
DEFAULT vs EDUCATION		
LIMIT_BAL vs SEX	Mann-Whitney U test and t test	P<0.05
LIMIT_BAL vs MARRIAGE		
LIMIT_BAL vs EDUCATION	Kruskai-Wallis test and ANOVA	P<0.05
LIMIT_BAL vs DEFAULT	Mann-Whitney U test and t test	P<0.05

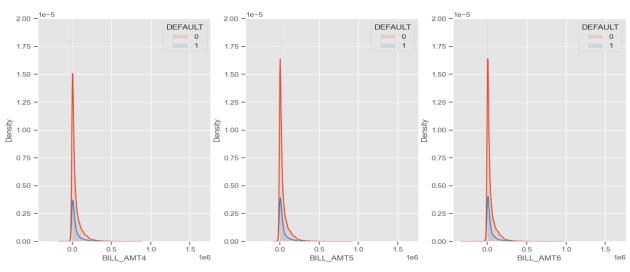
Data Exploration



Data Exploration







Feature Engineering

$$1.pay_sum = \sum_{i=1}^{6} PAY_i$$

$$2.mean_utilization_ratio = \sum_{i=1}^{6} BILL_AMTi / (6 \times LIMIT_BAL)$$

$$3.6_month_loss_given_default = \sum_{i=1}^{6} BILL_AMTi - \sum_{i=1}^{6} PAY_AMTi$$

$$4.mean_payment_ratio = \sum_{i=1}^{6} PAY_AMTi / (6 \times LIMIT_BAL)$$

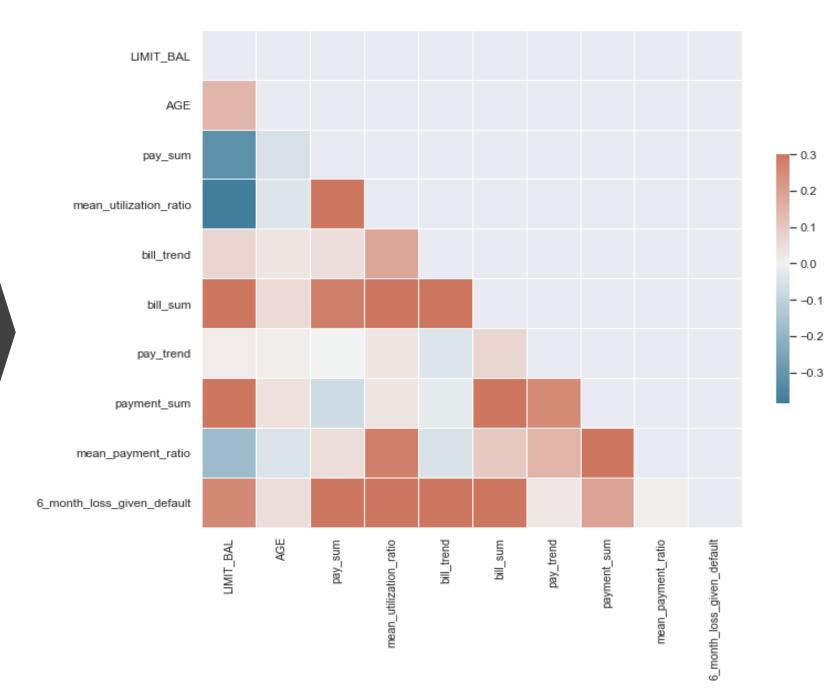
$$5.bill_trend = (\sum_{i=1}^{3} BILL_AMTi - \sum_{i=4}^{6} BILL_AMTi) / (3 \times LIMIT_BAL)$$

$$6.pay_trend = (\sum_{i=1}^{3} PAY_AMTi - \sum_{i=4}^{6} PAY_AMTi) / (3 \times LIMIT_BAL)$$

$$7.bill_sum = \sum_{i=1}^{6} BILL_AMTi$$

$$8.payment_sum = \sum_{i=1}^{6} PAY_AMTi$$

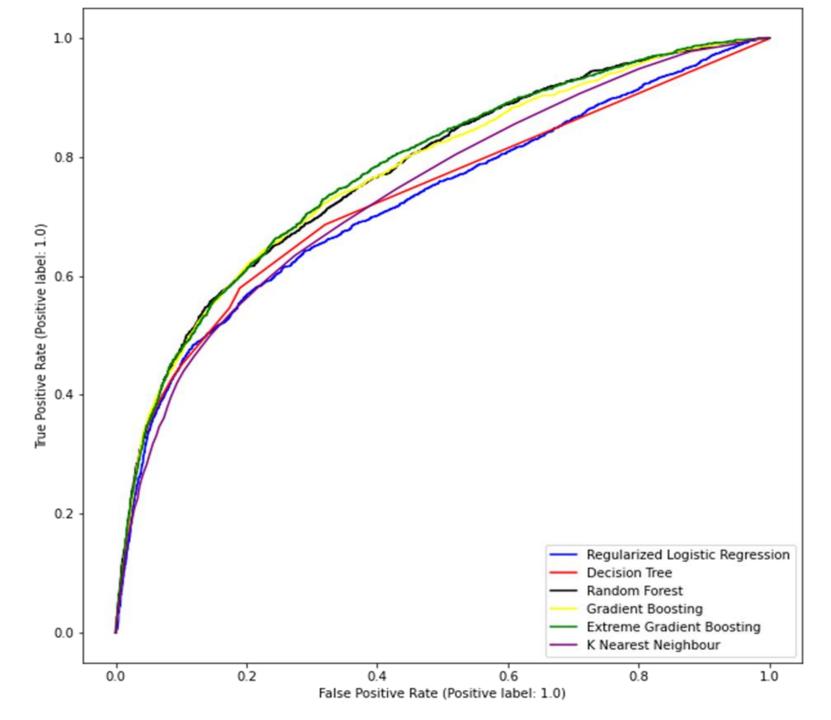
Feature Engineering



Modeling using the original data set

Algorithms	Hyperparameters	best estimate	accuracy	precision	recall	f1 score	AUC	runtime	top 3 most important features
regularized logistic regression	С	50	0.809	0.68	0.234	0.348	0.725	4.77	PAY_1, BILL_AMT1, BILL_AMT2
decision tree	criterion, max_depth	gini, 4	0.819	0.66	0.349	0.457	0.725	36.73	PAY_1, PAY_2, PAY_AMT3
random forest	n_estimators, max_depth	200, 9	0.821	0.67	0.356	0.465	0.774	85.31	PAY_1, PAY_2, PAY_3
gradient boosting	learning_rate, n_estimator, max_depth	50, 2, 0.1	0.822	0.68	0.349	0.46	0.771	1926.3	PAY_1, PAY_2, BILL_AMT1
Extreme gradient boosting	subsamples, n_estimators, max_depths, learning rate, gamma, reg_alpha	0.5, 200, 5, 0.05, 0.1, 0	0.819	0.66	0.358	0.463	0.778	301.9	PAY_1, PAY_1, PAY_3
KNN	n_neighbors, p	50, 2	0.808	0.646	0.258	0.368	0.741	100.8	PAY_1, PAY_2, PAY_3

ROC curves comparison



Modeling using engineered features

Algorithms	Hyper- parameters	best estimate	accuracy	precision	recall	f1 score	AUC	top 3 most important features
regularized logistic regression	C	50	0.797	0.672	0.176	0.279	0.692	pay_sum, payment_sum, LIMIT_BAL
decision tree	criterion, max_depth	entropy, 4	0.801	0.625	0.276	0.383	0.743	pay_sum, payment_sum, bill_sum
random forest	n_estimators, max_depth	200, 9	0.801	0.612	0.299	0.398	0.762	<pre>pay_sum, payment_sum, bill_trend</pre>
gradient boosting	learning_rate, n_estimator, max_depth	50, 2, 0.1	0.802	0.616	0.299	0.402	0.759	pay_sum, payment_sum, bill_sum
Extreme gradient boosting	subsamples, reg_alpha, n_estimators, max_depths, learning rate, gamma	0.7, 0.05, 400, 1, 0.1,0.4	0.803	0.614	0.322	0.422	0.759	pay_sum, payment_sum, bill_sum
KNN	n_neighbors, p	2, 50	0.796	0.636	0.204	0.309	0.713	pay_sum, payment_sum

Conclusion

- Regardless of the socioeconomic status of the clients, payment status is the most critical feature for credit default prediction. It is described as PAY_1 to PAY_6 in the original data or pay_sum in the reduced data set.
- Features related to a socioeconomic status like age, marriage, and education significantly affect the default in credit risk assessment.
- XGB and random forest are the most attractive algorithm for predicting credit default risk compared with RLR, DT, GB, and KNN.

Limitations

- Data bias: Because this data set is from Taiwan instead of the US, it has limited application reference for consumer credit prediction in the US.
- All these analysis is only applied to the existing clients of the credit card company, not for the prospective ones.
- Features: We didn't have credit bureau data in this project. We also applied historical data, not the information recently.

Future work



Considering the popularity of deep learning model, we can explore its application in this data set and compare it with the ensemble algorithms.



Develop a modeling pipeline to extract information more efficiently, providing an automated and faster solution for making credit decisions on time.