

## Credit Default Risk Analysis

Lina Gao

## Problem Statement

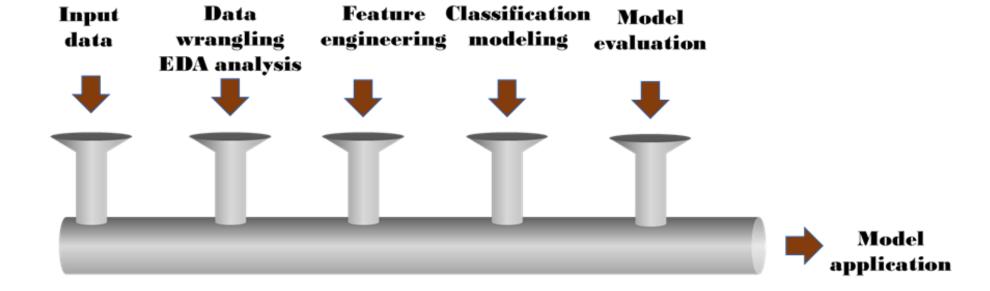
• How do we use modern machine learning algorithms to identify potential credit defaulter based on their historical transaction data and socioeconomic status?

#### Source of Data

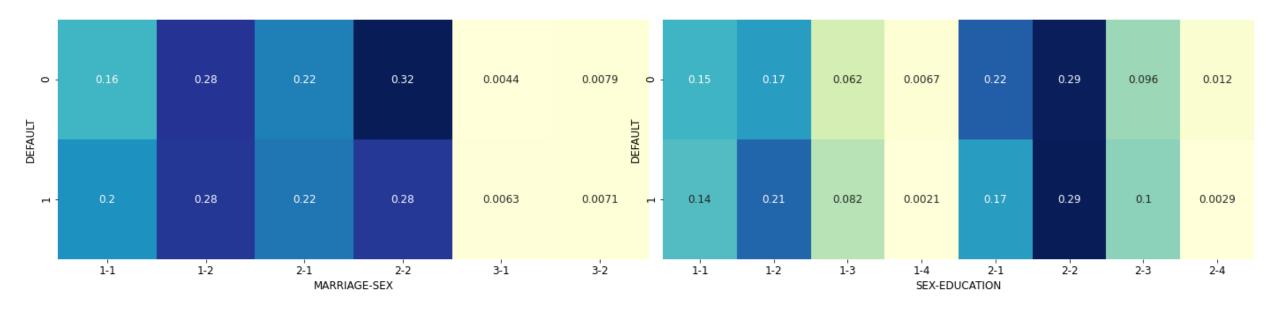
- ➤ It includes 30000 observations with 25 features collected in Taiwan from April 2005 to Sep 2005.
- Features: Age, Sex, Marriage, Education, credit limit, payment, bill and payment status in each of the six months.



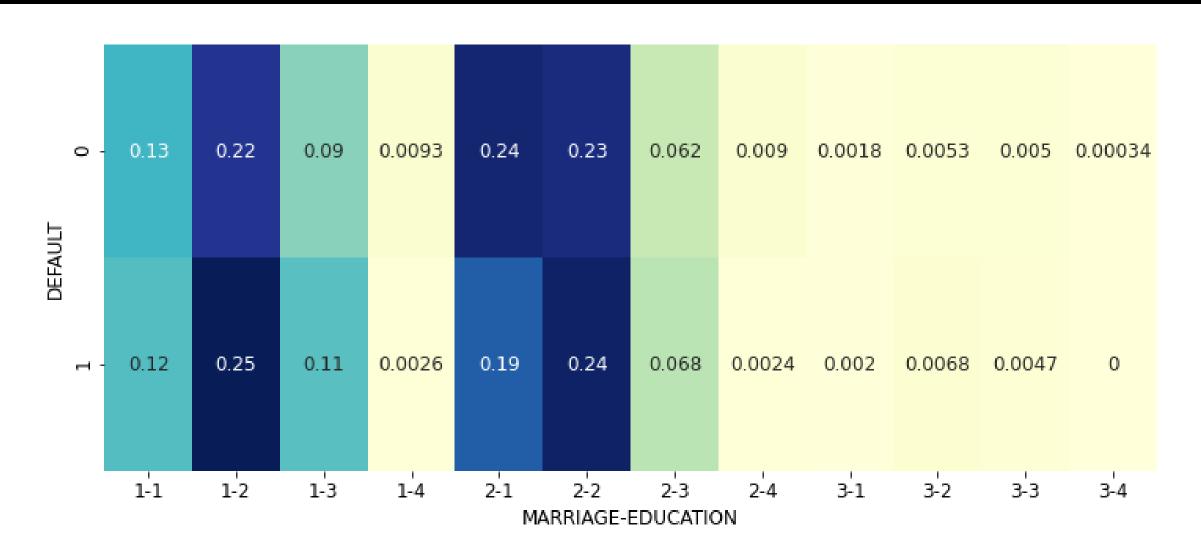
## Analytical Workflow



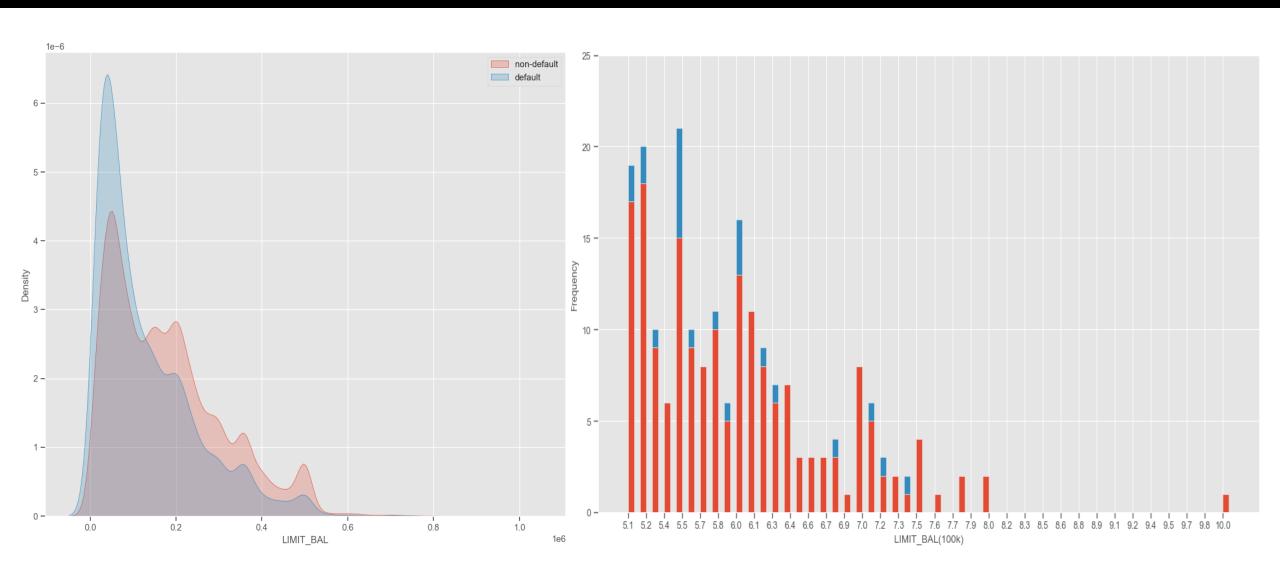
#### Data Exploration: client segmentation



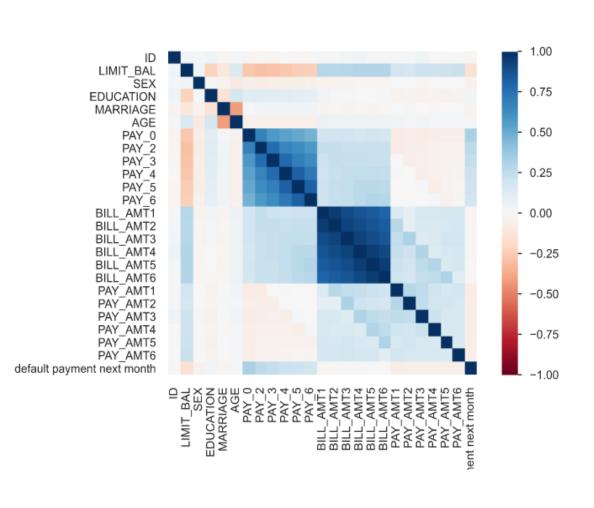
#### Data Exploration

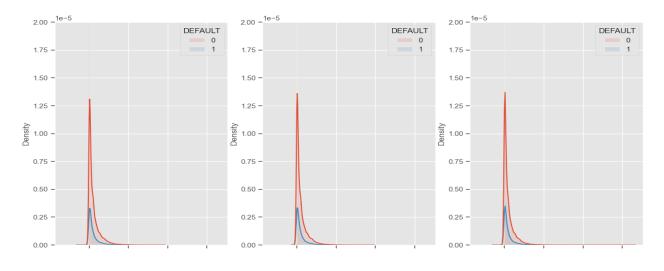


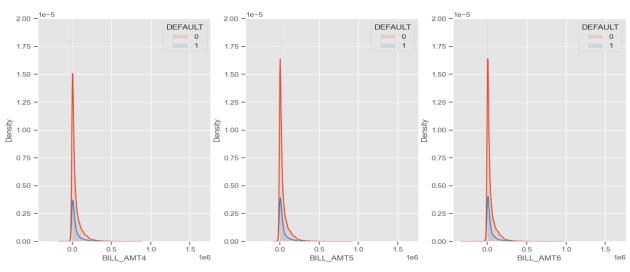
### Data Exploration



#### **Data Exploration**







## Feature Engineering

$$1.pay\_sum = \sum_{i=1}^{6} PAY\_i$$

$$2.mean\_utilization\_ratio = \sum_{i=1}^{6} BILL\_AMTi / (6 \times LIMIT\_BAL)$$

$$3.6\_month\_loss\_given\_default = \sum_{i=1}^{6} BILL\_AMTi - \sum_{i=1}^{6} PAY\_AMTi$$

$$4.mean\_payment\_ratio = \sum_{i=1}^{6} PAY\_AMTi / (6 \times LIMIT\_BAL)$$

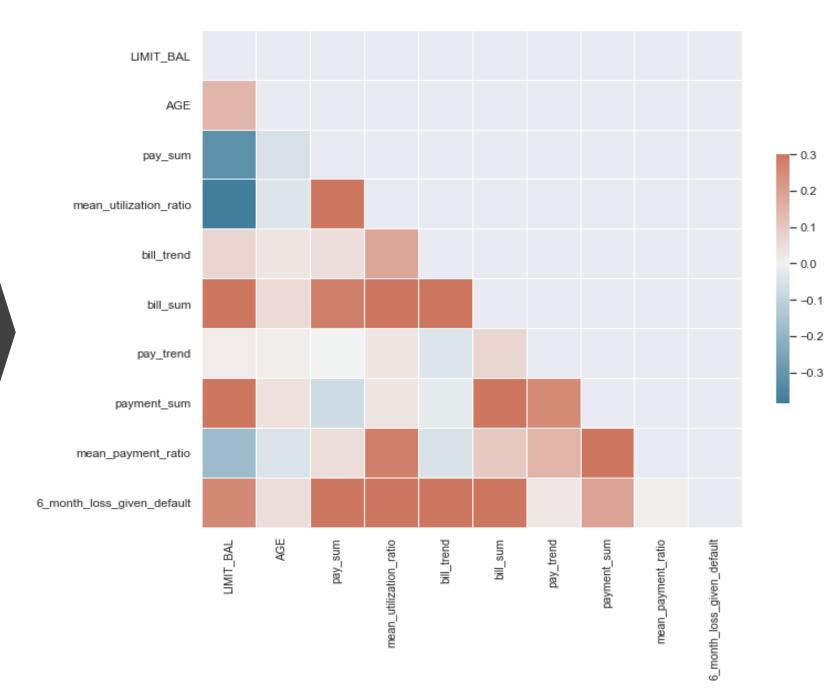
$$5.bill\_trend = (\sum_{i=1}^{3} BILL\_AMTi - \sum_{i=4}^{6} BILL\_AMTi ) / (3 \times LIMIT\_BAL)$$

$$6.pay\_trend = (\sum_{i=1}^{3} PAY\_AMTi - \sum_{i=4}^{6} PAY\_AMTi ) / (3 \times LIMIT\_BAL)$$

$$7.bill\_sum = \sum_{i=1}^{6} BILL\_AMTi$$

$$8.payment\_sum = \sum_{i=1}^{6} PAY\_AMTi$$

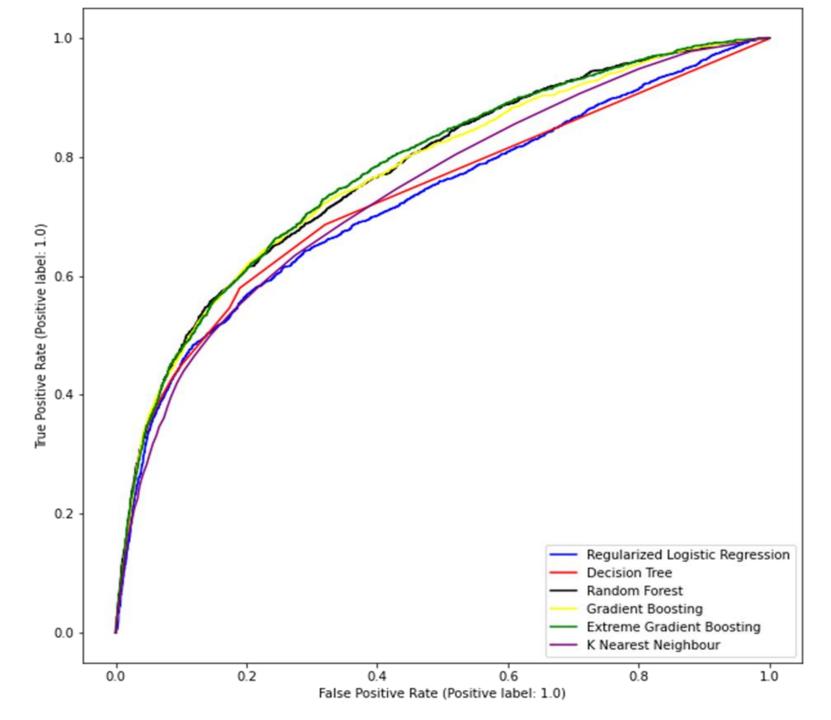
## Feature Engineering



## 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.