

# Generating Simulated Financial Data

Yitong Wang, Yishuang Chen, Xi Fan, Yuxin Zhang, Qi Chen

## PROJECT OVERVIEW

Together with Generative Adversarial Networks, we can say "No" to limited data, but say "Hi" to infinite supplementary data

## ENCODING

The original data set includes 1,103,307 records and 44 features. After data cleaning and feature selection, we got 11,936 records and 7 features. By label encoding, we transform all the continuous and categorical data into letters

ACTION	CLEARED	COLLATERALIZATION	TAXONOMY	TRADE CONTINUATION	PRICE	AMOUNT
0 NEW	U	UC	CreditIndex.CDX.CDXIG	Trade	1	100
1 NEW	U	PC	CreditIndex.CDX.CDXIG	Trade	1	30
2 CORRECT	U	PC	CreditIndex.CDX.CDXIG	Trade	1	30
3 NEW	U	PC	CreditIndex.CDX.CDXIG	Novation	1.24	25
4 CORRECT	C	UC	CreditIndex.CDX.CDXHY	Trade	1	5
5 CORRECT	C	FC	CreditIndex.CDX.CDXIG	Trade	1	58
NE			IG	T	(0,1]	a m : (90,100]

Sample of encoded text:

NE U UC IG T a m, NE U PC IG T a f, CR U PC IG T a f...

We build a vocabulary consisting of 847 combinations of the 7 features. Then we map each possible value into number. We aggregate every 30 orders as one sequence, finally got 385 sequences.

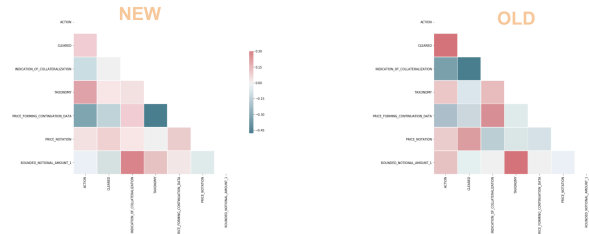
Action  
"NEW": "NE"  
"CORRECT": "CR"  
"CANCEL": "CC"  
denotes the action we operate on the CDS index

Collateralization  
"UC": Uncollateralized  
"FC": Fully collateralized  
"PC": Partially collateralized  
"OC": One-way collateralized  
reflects the risk of CDS index

Trade Continuation  
"Amendment": A  
"Novation": N & "Trade": T  
"Partial termination": P  
denotes the post-execution trade continuation or life cycle

## EVALUATION METRIC

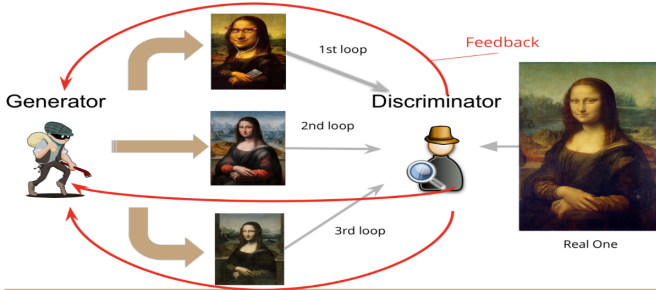
To evaluate the quality of generated data, we need to design a metric. The first measurement is to check whether the correlation among different features still holds. We plot a correlation matrix between each pair of features.



The results show that the low correlation between different feature still holds

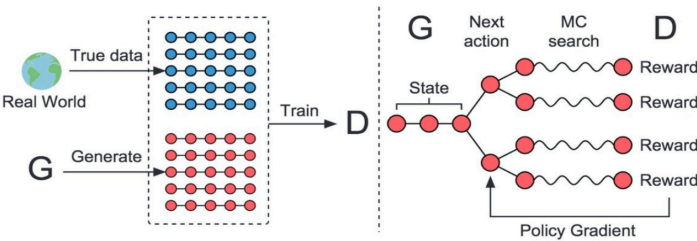
## MODEL

Generative Adversarial Network (GAN) has succeed in generating synthetic pictures that mimic the real one. We want to use this advantage to generate more financial data which has the same attribute as the real data to solve the data lacking problem. In that way, we can use the generated data for better financial model building.



## MODEL -- SeqGAN

We process our financial data as sequence data. However, the vanilla GAN has problems in generating discrete data. We applied the seqGAN built by Lantao Yu to remodel our generator and discriminator in a reinforcement learning framework.



(<https://arxiv.org/abs/1609.05473?context=cs.AI>)

Next we want to see whether the relative performance of each two algorithms (trained and tested) on the synthetic dataset is the same as their relative performance (when trained and tested) on the original dataset. We use five different models to classify the "cleared" feature. And then we use a metric called SRA to evaluate that.

Accuracy	Random Forest	Logistic Regression	Linear SVM	Multinomial Naive Bayes	Gaussian Naive Bayes
Original Data	0.87144	0.74539	0.69305	0.62353	0.57789
Generated Data	0.983	0.79265	0.7886	0.64165	0.3268

$$R_i = m_{\tau}(A_i(D_1), D_2)$$
$$S_i = m_{\tau}(A_i(D_1^C), D_2^C)$$

$R_i, S_i$  Performance metric score  
 $m_{\tau}$  Performance metric on task  $\tau$   
 $A_i(D_1)$  Trained model from training data  
 $D_1^C$  Training data from synthetic data  
 $D_2^C$  Test data from synthetic data

$SRA(G) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n I((R_i - R_j) \times (S_i - S_j) > 0)$   
The SRA can be thought of as the (empirical) probability of a comparison on the synthetic data being "correct".  $SRA(\text{"Cleared"}) = 1$

Other features:

$SRA(\text{"Action"}) = 0.8$ ,  $SRA(\text{"Collateralization"}) = 0.9$ ,  $SRA(\text{"Price forming"}) = 0.9$ ,  $SRA(\text{"Taxonomy"}) = 0.6$   
Based on the result, we can see the quality of the synthetic data is good in total.

## Future work

Model | Hyperparameter tuning  
Features | Underlying asset, execution venue, design new features  
Application | Strategic design, simulation

## Acknowledgement

J.P.Morgan

FUNG INSTITUTE FOR  
ENGINEERING LEADERSHIP  
UC BERKELEY ENGINEERING

BERKELEY  
IEOR  
INDUSTRIAL ENGINEERING  
& OPERATIONS RESEARCH



Guo Xin | Professor at UCB

Cheng-Ju Wu, Yang Nan | Ph.D student at UCB

Ian Goodfellow | Father of GAN

Lantao Yu | Primary author of SeqGAN paper