

# STAT 430 Deep Learning Final project

1<sup>st</sup> Yuxuan Wan

*Department of Statistics*

*University of Illinois Urbana-Champaign*

Champaign, USA

yuxuanw8@illinois.edu

2<sup>nd</sup> Min-Hsueh Chiang

*Department of Statistics*

*University of Illinois Urbana-Champaign*

Champaign, USA

mchiang5@illinois.edu

3<sup>rd</sup> Xiao Ma

*Department of Statistics*

*University of Illinois Urbana-Champaign*

Champaign, USA

xm22@illinois.edu

Team ID: 15

Team Name : transfer-learning

Team Member Names : Yuxuan Wan, Min-Hsueh Chiang, Xiao Ma

Team member NetIDs : yuxuanw8, mchiang5, xm22

**Abstract**—This document is a summary of the project of credit card fraud detection. Main idea is applying the deep learning techniques to the detection of credit card fraud. Our method is to build a n-layers deep neural network to improve the accuracy of credit card fraud. Thus, one of the most important ideas is deciding how many layers we need in our model. Another important point is how we compile and train the model. The implications we have done so far is building the foundation for the future work. Based on the method of DNN, we need to figure out how many layers we will set. And, we need to compile and train the mode in future work.

## I. INTRODUCTION

Our problem being solved is the detection of fraudulent transactions in credit card payments. With the fast-processing internet, online shopping has become a more fraudulent transaction common way of shopping nowadays. Moreover, a credit card is the most common way of paying online, but there are still some drawbacks to shopping online. For instance, fraudulent transactions in online systems are constantly increasing as well, especially for credit cards.

Following on, Credit card fraud detection is important for both the user end and the bank or credit card company end. For the user end, more accurate credit card fraud detection can result in better credit card customer experiences and help protect customers from credit card fraud which leads to money loss. For the credit card company, using deep learning techniques to improve the accuracy of recognizing and detecting credit card fraud can protect companies' assets as well as companies' reputations so that they can provide customers with better credit card service and user experiences.

Based on our method, the input(predictor variable) of our project is the dataset created by Banksim which is a simulator of bank payment based on a sample of aggregated bank data. The data set contains 594643 observations and 10 columns. The output is generating the label(Mostly, binary variable) of fraudulent transactions. And, this is to differentiate and classify fraudulent transactions and achieve our goal(identifying the fraudulent transactions).

## II. DATA

We got access to the data set from Kaggle, provided by EDGAR LOPEZ-ROJAS. However, the origin of the data was created in the 26Th European Modeling and Simulation Symposium, EMSS 2014, Bordeaux, France, pp. 144–152, Dime University of Genoa, 2014, ISBN: 9788897999324. The data set was created by Bank-sim which is a simulator of bank payment based on a sample of aggregated bank data. Moreover, the main purpose of generating synthetic data from the simulator is to detect potential fraudulent transactions happening nowadays. The data set contains 594643 observations with 10 columns, features, in total. As we can see that the step, amount, and the fraud column are the only columns that are numeric in this data set, whereas the other remaining columns are all object type columns.

For the data-pre-processing part, we have printed the data and discovered that there are some unknown values in both the age and the gender columns. To be specific, there is a character "U" existing in both columns, which does not seem to be normal to both the age and the gender column. Thus, we decided to replace the U with age 2 and gender Male, in order to make the data more validated. Moreover, we decided to transform the age column from an object type of data to a numeric type of data, in which it makes more sense for the age being a numerical variable. In addition, we also subset the data set into two different data frames in which one contains only the fraud data and the other contains only the non-fraud data. Following on, we provided a brief view of what the data set looks like, the first five rows of the data, a frequency table which demonstrates there are 587443 of non-fraud observations and 7200 of fraud observations in total. Subsequently, we also provide a correlation matrix for the numerical predictor variables in order to check if there are multi-collinearity issues.

step	customer	age	gender	zipcodeOri	merchant	zipMerchant	category	amount	fraud
0	0	'C1093826151'	4	'M'	'28007'	'M348934600'	'28007' 'es_transportation'	4.55	0
1	0	'C352968107'	2	'M'	'28007'	'M348934600'	'28007' 'es_transportation'	39.68	0
2	0	'C2054744914'	4	'F'	'28007'	'M1823072687'	'28007' 'es_transportation'	26.89	0
3	0	'C1760612790'	3	'M'	'28007'	'M348934600'	'28007' 'es_transportation'	17.25	0
4	0	'C757503768'	5	'M'	'28007'	'M348934600'	'28007' 'es_transportation'	35.72	0

Fig. 1. First five observations of the data set.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 594643 entries, 0 to 594642
Data columns (total 10 columns):
#   Column      Non-Null Count  Dtype
---  -
0   step        594643 non-null  int64
1   customer    594643 non-null  object
2   age         594643 non-null  int64
3   gender      594643 non-null  object
4   zipcodeOri  594643 non-null  object
5   merchant    594643 non-null  object
6   zipMerchant 594643 non-null  object
7   category    594643 non-null  object
8   amount      594643 non-null  float64
9   fraud       594643 non-null  int64
dtypes: float64(1), int64(3), object(6)
memory usage: 45.4+ MB
```

Fig. 2. Types of data.

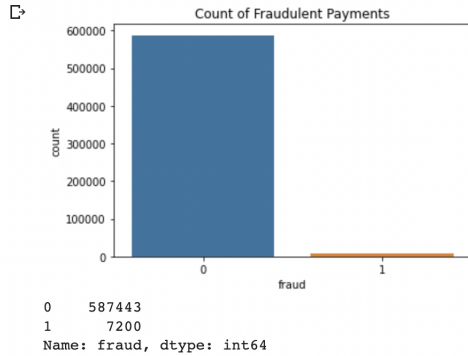


Fig. 3. Total fraud and non-fraud counts.

	age	amount	step
age	1.000000	-0.003878	0.000563
amount	-0.003878	1.000000	-0.007961
step	0.000563	-0.007961	1.000000

Fig. 4. Correlation Matrix.

### III. PRELIMINARY TECHNICAL DETAILS AND RESULTS

We are planning to apply a Dense Neural network (DNN) model in training our model and predicting for the testing dataset. To be specific, DNN is a deep learning technique which employs multiple layers and activations between the input and output layers, in which each layer is dependent to the previous layer. When implementing DNN, we have to first split our datasets into training and testing datasets, which the input vector, including the predictor variables, and our target which is the potential response variable should all be included as well. After splitting, we should define our neural network and the initialization stage in our model, in which we should include potential fully connected layers, and activation functions. We decided to elect the ReLU activation function for our overall neural network, because ReLU generally performs better than other activation functions in the DNN models, because it will not be easily affected by extreme values. Moreover, in order for the loss to be positive in our case, and to prevent overflow and underflow issues would employ a log softmax activation function before the output layer, in order to enhance the loss's positivity and validity.

Following on, we should load our training dataset into our training wrapper. During the training wrapper, we first initialize our data sets, and for each batch and each epochs, we employ a and Adam gradient descent technique which continuously run to minimize the error function until it converges to a predetermined lowest value, starting from an initialization stage where the model parameters are initialized to an initial set of values. Subsequently, we calculate the loss for every epoch with our chosen binary cross-entropy loss formula. Eventually, we load our testing data set to our trained model, and calculate for its average loss and model accuracy to verify if our model is a good one or not.

### IV. RELATED WORK

This part listed scholarly works that relate to the project and our views about these scholarly works:

A. Dal Pozzolo, Andrea. "Adaptive machine learning for credit card fraud detection." (2015). [1]

Dal Pozzolo's paper is trying to use techniques from machine learning, more specifically, supervised learning methods, to perform the classification problem [1]. It is aiming to use methodology from classification problems to classify data as credit card fraud detected transactions and credit card fraud-not-detected transactions by using the training data. Then, the cost-based measures are defined to assess the performance of the model constructed through training data and applied on the test dataset [1]. Our project is also trying to tackle the challenge of credit card fraud detection but through the deep learning method. The procedures of constructing the model and assessing the model constructed of Dal Pozzolo's project and our project are similar, but the methodology of constructing the model are different. We are trying to build the dense neural network (DNN) through the training data splitted from the total dataset as our model and then apply

that DNN model to the test dataset to assess the performance on the cost-based measure.

*B. Carcillo, Fabrizio, et al. "Combining unsupervised and supervised learning in credit card fraud detection." Information sciences 557 (2021): 317-331. [2]*

Fabrizio Carcillo's paper is trying to combine both unsupervised learning methodology and supervised learning methodology together to perform the credit card fraud detection [2]. Compared with Dal Pozzolo's approach [1], Fabrizio is going to apply supervised techniques which can learn from the past fraudulent transactions and behaviors to complement with unsupervised techniques which can assist to detect new types of fraud [2]. The high-level procedure methodology of Fabrizio's paper is similar to Dal Pozzolo's paper and ours. Splitting the dataset into training data and testing data, they use the training data portion to train the model and use the test dataset to assess the performance of the model on the cost-based measure which is also similar to what Dal Pozzolo proposed in his paper. It is the same story here. Fabrizio Carcillo is trying to solve the similar problem about credit card fraud detection but using different methods of constructing models and assessing the model [2].

*C. Lebichot, Bertrand, et al. "Deep-learning domain adaptation techniques for credit cards fraud detection." INNS Big Data and Deep Learning conference. Springer, Cham, 2019. [3]*

Bertrand Lebichot proposed using the combination of transfer learning methods, which can store the knowledge gained while solving one problem and using that knowledge gained and applying it to solve another problem related to the known problem, and domain adaptation [3]. Bertrand Lebichot discussed five methods BNNN, NDNN, FEDADNN, AguDNN, and AdvDnn and compared them on the credit card transaction dataset [3]. Then, he used precision@100 to compare the performance of those five different methods. This work relates to our project in the way that it is trying to apply the transfer learning methods and domain adaptation for which dense neural networks methods have been widely used [3]. We are also trying to utilize DNN to train our model so that we can better detect credit card fraud.

*D. Maes, Sam, et al. "Credit card fraud detection using Bayesian and neural networks." Proceedings of the 1st international nauso congress on neuro fuzzy technologies. Vol. 261. 2002. [4]*

This paper is trying to use two types of network to detect credit card fraud. One is Feed Forward Multi-layer Perceptron which uses error correction learning, and the other is the Bayesian network which is a probabilistic graphical model for representing knowledge about an uncertain domain where each node corresponds to a random variable and each edge represents the conditional probability for the corresponding random variable [4]. This paper relates to our project in the way that it discusses how Feed Forward Multi-layer Perceptron

contributes to the detection of credit card fraud. This is similar to the model we are going to construct to train our own neural network so that we can better identify credit card transactions as fraudulent.

*E. Roy, Abhimanyu, et al. "Deep learning detecting fraud in credit card transactions." 2018 Systems and Information Engineering Design Symposium (SIEDS). IEEE, 2018. [5]*

Abhimanyu Roy evaluated a subset of Deep Learning topologies with regard to their effectiveness in detecting credit card fraud on a dataset of nearly 80 million credit card transactions that have been pre-labeled as fraudulent and legitimate [5]. These topologies range from the general artificial intelligence neural networks to topologies with built-in time and memory components, such as Long Short-term memory [5]. To get beyond typical fraud detection issues like class imbalance and scalability, they used a high speed, distributed cloud computing platform [5]. Their study offers a thorough manual for model parameter sensitivity analysis with respect to credit card fraud detection performance. In order to help financial institutions cut losses by as much as 50%, they also propose a methodology for parameter tuning Deep Learning topologies for credit card fraud detection [5]. This paper relates to our project in the way that it also utilizes Feed Forward Multilayer Perceptron neural network which consists of different layers of perceptrons that are interconnected by a set of connections which have different weights associated with them.

*F. Pumsirirat, Apapan, and Yan Liu. "Credit card fraud detection using deep learning based on auto-encoder and restricted boltzmann machine." International Journal of advanced computer science and applications 9.1 (2018). [6]*

Apapan Pumsirirat and Yan Liu, in their paper, are trying to use deep learning which is based on auto-encoder and restricted Boltzmann Machine to detect fraud cases that cannot be detected based on previous history or supervised learning [6]. They built a deep auto-encoder and restricted Boltzmann machine (RBM) model that can reassemble typical transactions and look for anomalies in the course of regular patterns [6]. This paper is helpful for our project topic, since it discusses how supervised learning, especially like the multi-layer perceptron model, fails to capture the type of credit card that has never happened before, since supervised learning is based on the previous history. This is the point that we should consider in order to empower fraud detection system, since only detecting prior fraud is not enough.

## V. CONTRIBUTION

Xiao Ma wrote out the abstract and introduction part of this progress report. Also, he wrote out the references. Min-Hsueh Chiang wrote out the data part of the progress report. Moreover, he wrote out the preliminary technical details. Yuxuan Wan wrote out the related work and how they relate to our project. Everyone is contributing equally at this proposal stage of the project. Contribution percentage: each person 33.333

## VI. APPENDIX

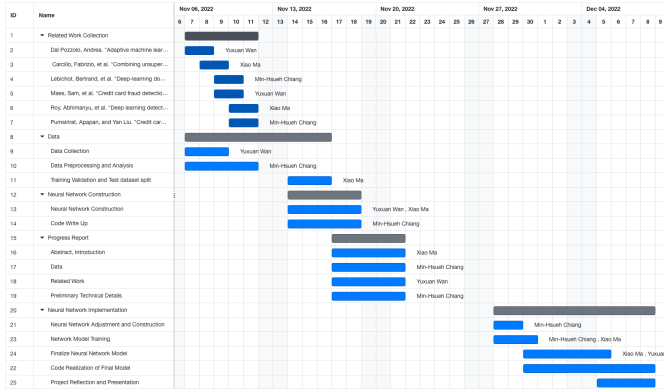


Fig. 5. Timeline of Work

## REFERENCES

- [1] Dal Pozzolo, A. (2015). Adaptive machine learning for credit card fraud detection.
- [2] Carcillo, F., Le Borgne, Y. A., Caelen, O., Kessaci, Y., Oblé, F., & Bontempi, G. (2021). Combining unsupervised and supervised learning in credit card fraud detection. *Information sciences*, 557, 317-331.
- [3] Lebichot, B., Borgne, Y. A. L., He-Guelton, L., Oblé, F., & Bontempi, G. (2019, April). Deep-learning domain adaptation techniques for credit cards fraud detection. In *INNS Big Data and Deep Learning conference* (pp. 78-88). Springer, Cham.
- [4] Maes, S., Tuyts, K., Vanschoenwinkel, B., & Manderick, B. (2002, January). Credit card fraud detection using Bayesian and neural networks. In *Proceedings of the 1st international naisto congress on neuro fuzzy technologies* (Vol. 261, p. 270).
- [5] Roy, A., Sun, J., Mahoney, R., Alonzi, L., Adams, S., & Beling, P. (2018, April). Deep learning detecting fraud in credit card transactions. In *2018 Systems and Information Engineering Design Symposium (SIEDS)* (pp. 129-134). IEEE.
- [6] Pumsirirat, A., & Liu, Y. (2018). Credit card fraud detection using deep learning based on auto-encoder and restricted boltzmann machine. *International Journal of advanced computer science and applications*, 9(1).