

Financial Fraud Detection

Subjects: [Computer Science, Information Systems](#)

Contributor: Patience Chew Yee Cheah , Yue Yang , Boon Giin Lee

The financial sector faces a significant challenge in the form of financial fraud, encompassing various forms of criminal deception aimed at securing financial gains, including activities like telecommunication fraud and credit card skimming. The proliferation of electronic payment technology has propelled online transactions into the mainstream, thereby amplifying the occurrence of fraudulent schemes. The prevalence of these fraudulent transactions has led to substantial losses for financial institutions. However, the large daily transactions pose a challenge for humans in manually identifying fraud.

[financial fraud detection](#)

[class imbalance](#)

[SMOTE](#)

[GANs](#)

[SMOTified-GAN](#)

1. Introduction

Recently, deep learning techniques have been explored and have shown promising results in detecting financial fraud [Alarfaj et al. \(2022\)](#); [Fang et al. \(2021\)](#); [Kim et al. \(2019\)](#). Unfortunately, most real-world financial fraud datasets suffer from a severe class imbalance issue, where the fraud data's proportion is significantly lower than that of nonfraud. In binary classification, class imbalance often leads to biased predictions favoring the majority class [Johnson and Khoshgoftaar \(2019\)](#). Consequently, the classifier's performance on the minority class is compromised, especially when encountering dissimilar frauds. Overcoming this problem poses a significant challenge, as classifiers are expected to achieve high precision and recall in fraudulent class.

To address this problem, several oversampling methods have been employed to generate minority samples. Synthetic Minority Oversampling TEchnique (SMOTE) interpolates between the existing minority data to synthesize minority samples [Chawla et al. \(2002\)](#). Generative Adversarial Networks (GANs) comprise a discriminator that aims to differentiate between real and generated samples and a generator that strives to deceive the discriminator by synthesizing realistic samples [Goodfellow et al. \(2014\)](#). GANs have shown superior results compared with SMOTE [Fiore et al. \(2019\)](#). However, SMOTE may cause overgeneralization issues. GAN, primarily designed for image generation, is not ideal for handling the class imbalance problem. To overcome these limitations, SMOTified-GAN employs SMOTE-generated samples instead of random noises as input to the GAN [Sharma et al. \(2022\)](#).

2. Detecting Method for Financial Fraud

The task of detecting financial fraud can be approached as a binary classification challenge, where classifiers examine the patterns within fraudulent and legitimate transactions to classify new transactions accurately. Consequently, it is crucial to possess an ample and diverse dataset to enable classifiers to grasp the inherent

patterns of both transaction categories. Addressing the issue of inadequate fraudulent samples in the training dataset, various methodologies have been introduced to create artificial fraud instances and supplement the original data. These techniques include SMOTE, GAN, and SMOTified-GAN.

SMOTE [Chawla et al. \(2002\)](#) has been widely applied to imbalanced training datasets. More than 85 SMOTE variations were proposed by 2018, including SMOTE+TomekLinks, SMOTE+ENN, Borderline-SMOTE, and Adaptive Synthetic [Fernández et al. \(2018\)](#). Recent studies proposed Radius-SMOTE [Pradipta et al. \(2021\)](#), which prevents overlap among generated samples, and Reduced-Noise SMOTE [Arafa et al. \(2022\)](#), which removes noise after oversampling. In financial fraud detection, SMOTE and its variations have been widely utilized to resample highly imbalanced datasets before training models such as AdaBoost [Ileberi et al. \(2021\)](#) and FNN [Fang et al. \(2021\)](#). Besides the finance domain, SMOTE and its variations have found extensive application in other fields dealing with highly imbalanced datasets. In bio-informatics, SMOTE has been used to discriminate Golgi proteins [Tahir et al. \(2020\)](#) and predict binding hot spots in protein–RNA interactions [Zhou et al. \(2022\)](#). In medical diagnosis, SMOTE and its variations have been employed for diagnosing cervical cancer [Abdoh et al. \(2018\)](#) and prostate cancer [Abraham and Nair \(2018\)](#). SMOTE has also been used to predict diabetes [Mirza et al. \(2018\)](#) and heart failure patients' survival [Ishaq et al. \(2021\)](#).

GANs [Goodfellow et al. \(2014\)](#) and their variations have more recently been employed for generating minority samples to tackle the class imbalance problem. [Douzas and Bacao \(2018\)](#) utilized a conditional GAN (cGAN) which can recover the distribution of training data to generate minority samples. To address the mode collapse issue, Balancing GAN was proposed to generate more diverse and higher-quality minority images [Mariani et al. \(2018\)](#). However, in this technique, the generator and discriminator cannot simultaneously reach their optimal states, leading to the development of IDA-GAN [Yang and Zhou \(2021\)](#). In financial fraud detection, GAN has been employed to generate fraud samples for imbalanced datasets before training classifiers, such as AdaBoost-Decision Tree [Mo et al. \(2019\)](#) and FNN [Fiore et al. \(2019\)](#). These studies have reported that the GAN achieves higher AUC, accuracy, and precision compared with SMOTE. Interestingly, [Fiore et al. \(2019\)](#) found that the best performance was achieved when twice as many GAN-generated fraud samples as the original fraud data were added to the training dataset. In other finance-related domains, GANs have been utilized to address class imbalance in money laundering detection in gambling [Charitou et al. \(2021\)](#). GANs and their variations have also been used extensively for high-dimensional imbalanced datasets, such as images [Mariani et al. \(2018\)](#); [Scott and Plested \(2019\)](#) and biomedical data [Zhang et al. \(2018\)](#). Recent studies have successfully applied GANs and their variations to generate minority samples in bio-informatics [Lan et al. \(2020\)](#).

Despite the notable accomplishments of SMOTE and GAN, these methods have certain limitations. SMOTE may introduce noise that leads to overgeneralization [Bunkhumpornpat et al. \(2009\)](#). While GANs can generate more “realistic” data, they may not be ideal for handling imbalanced data, as it was originally designed for generating images using random noise. Additionally, there may be insufficient real minority data available for training the GAN [Mariani et al. \(2018\)](#). To address these limitations, [Sharma et al. \(2022\)](#) proposed SMOTified-GAN, which employs SMOTE-generated samples as input for GAN instead of random numbers, resulting in improved performance compared with SMOTE and GAN.

In early studies, financial fraud detection systems predominantly depended on rule-based methodologies, wherein human expertise in fraud was translated into rules to anticipate fraudulent activities [Zhu et al. \(2021\)](#). However, the evolving behaviors of fraudsters and the increasing size of transaction datasets have posed challenges in identifying fraud-related rules manually. As a result, research has shifted towards machine learning methods, such as naive Bayes, logistic regression, support vector machine, random forest, and decision tree ([Ileberi et al. 2021](#); [Ye et al. 2019](#); [Zhu et al. 2021](#)), which can “learn” fraud and nonfraud patterns from given datasets. Nonetheless, machine learning techniques require extensive data preprocessing before training the classifier [Alarfaj et al. \(2022\)](#); [Kim et al. \(2019\)](#); [Zhu et al. \(2021\)](#).

In recent years, deep learning has gained popularity in financial fraud detection due to its superior performance compared with traditional machine learning approaches [Alarfaj et al. \(2022\)](#); [Fang et al. \(2021\)](#); [Jurgovsky et al. \(2018\)](#); [Kim et al. \(2019\)](#). Some studies have approached financial fraud detection as a sequence classification problem, considering the temporal sequence of transactions as a crucial factor. Sequential models, such as Gated Recurrent Units [Branco et al. \(2020\)](#), Long Short-Term Memory (LSTM) [Jurgovsky et al. \(2018\)](#), and Time-aware Attention-based Interactive LSTM [Xie et al. \(2022\)](#), have been proposed. However, since most available financial fraud datasets lack time-sequence information, sequential models may not be suitable in such cases. Due to the vector format of finance fraud datasets without time-sequence information, FNNs are considered a suitable choice [Fang et al. \(2021\)](#); [Fiore et al. \(2019\)](#); [Kim et al. \(2019\)](#). Initially designed for image processing and classification, CNNs have also been found effective in financial fraud detection [Alarfaj et al. \(2022\)](#); [Chen and Lai \(2021\)](#); [Zhang et al. \(2018\)](#). Their 1D convolution layers can extract patterns within smaller segments of a transaction vector.

Building on [Fiore et al. \(2019\)](#)’s findings, this research aimed to assess the performance of a model using varying amounts of minority samples in the training dataset. To achieve this, the study explores the use of SMOTE, GAN, SMOTified-GAN, and other variants of hybrid SMOTE and GAN. Consequently, a combination of SMOTE- and GAN-generated minority samples, along with GANified-SMOTE, was proposed to fulfill the research aims. Finally, FNN, CNN, and FNN+CNN models were employed to ensure a fair evaluation of the performances of different data generation techniques.

References

1. Alarfaj, Fawaz Khaled, Iqra Malik, Hikmat Ullah Khan, Naif Almusallam, Muhammad Ramzan, and Muzamil Ahmed. 2022. Credit card fraud detection using state-of-the-art machine learning and deep learning algorithms. *IEEE Access* 10: 39700–15.
2. Fang, Weiwei, Xin Li, Ping Zhou, Jingwen Yan, Dazhi Jiang, and Teng Zhou. 2021. Deep learning anti-fraud model for internet loan: Where we are going. *IEEE Access* 9: 9777–84.
3. Kim, Eunji, Jehyuk Lee, Hunsik Shin, Hoseong Yang, Sungzoon Cho, Seung kwan Nam, Youngmi Song, Jeong a Yoon, and Jong il Kim. 2019. Champion-challenger analysis for credit card fraud detection: Hybrid ensemble and deep learning. *Expert Systems with Applications* 128: 214–24.

4. Johnson, Justin M., and Taghi M. Khoshgoftaar. 2019. Survey on deep learning with class imbalance. *Journal of Big Data* 6: 27.
5. Chawla, Nitesh V., Kevin W. Bowyer, Lawrence O. Hall, and W. Philip Kegelmeyer. 2002. SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research* 16: 321–57.
6. Goodfellow, Ian, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. *Advances in Neural Information Processing Systems* 27: 2672–80.
7. Fiore, Ugo, Alfredo De Santis, Francesca Perla, Paolo Zanetti, and Francesco Palmieri. 2019. Using generative adversarial networks for improving classification effectiveness in credit card fraud detection. *Information Sciences* 479: 448–55.
8. Sharma, Anuraganand, Prabhat Kumar Singh, and Rohitash Chandra. 2022. SMOTified-GAN for class imbalanced pattern classification problems. *IEEE Access* 10: 30655–65.
9. Fernández, Alberto, Salvador Garcia, Francisco Herrera, and Nitesh V. Chawla. 2018. SMOTE for learning from imbalanced data: Progress and challenges, marking the 15-year anniversary. *Journal of Artificial Intelligence Research* 61: 863–905.
10. Pradipta, Gede Angga, Retantyo Wardoyo, Aina Musdholifah, and I. Nyoman Hariyasa Sanjaya. 2021. Radius-SMOTE: A new oversampling technique of minority samples based on radius distance for learning from imbalanced data. *IEEE Access* 9: 74763–77.
11. Arafa, Ahmed, Nawal El-Fishawy, Mohammed Badawy, and Marwa Radad. 2022. RN-SMOTE: Reduced noise SMOTE based on DBSCAN for enhancing imbalanced data classification. *Journal of King Saud University—Computer and Information Sciences* 34: 5059–74.
12. Ileberi, Emmanuel, Yanxia Sun, and Zenghui Wang. 2021. Performance evaluation of machine learning methods for credit card fraud detection using SMOTE and AdaBoost. *IEEE Access* 9: 165286–94.
13. Tahir, Muhammad, Fazlullah Khan, Mohammad Khalid Imam Rahmani, and Vinh Truong Hoang. 2020. Discrimination of golgi proteins through efficient exploitation of hybrid feature spaces coupled with SMOTE and ensemble of support vector machine. *IEEE Access* 8: 206028–38.
14. Zhou, Tong, Jie Rong, Yang Liu, Weikang Gong, and Chunhua Li. 2022. An ensemble approach to predict binding hotspots in protein–RNA interactions based on SMOTE data balancing and random grouping feature selection strategies. *Bioinformatics* 38: 2452–58.
15. Abdo, Sherif F., Mohamed Abo Rizka, and Fahima A. Maghraby. 2018. Cervical cancer diagnosis using random forest classifier with SMOTE and feature reduction techniques. *IEEE Access* 6: 59475–85.

16. Abraham, Bejoy, and Madhu S. Nair. 2018. Computer-aided diagnosis of clinically significant prostate cancer from MRI images using sparse autoencoder and random forest classifier. *Biocybernetics and Biomedical Engineering* 38: 733–44.
17. Mirza, Shuja, Sonu Mittal, and Majid Zaman. 2018. Decision support predictive model for prognosis of diabetes using SMOTE and decision tree. *International Journal of Applied Engineering Research* 13: 9277–82.
18. Ishaq, Abid, Saima Sadiq, Muhammad Umer, Saleem Ullah, Seyedali Mirjalili, Vaibhav Rupapara, and Michele Nappi. 2021. Improving the prediction of heart failure patients' survival using SMOTE and effective data mining techniques. *IEEE Access* 9: 39707–16.
19. Douzas, Georgios, and Fernando Bacao. 2018. Effective data generation for imbalanced learning using conditional generative adversarial networks. *Expert Systems with Applications* 91: 464–71.
20. Mariani, Giovanni, Florian Scheidegger, Roxana Istrate, Costas Bekas, and Cristiano Malossi. 2018. BAGAN: Data augmentation with balancing GAN. *arXiv arXiv:1803.09655*.
21. Yang, Hao, and Yun Zhou. 2021. IDA-GAN: A novel imbalanced data augmentation GAN. Paper presented at the 2020 25th International Conference on Pattern Recognition (ICPR), Milan, Italy, January 10–15; pp. 8299–305.
22. Mo, Zan, Yanrong Gai, and Guanlong Fan. 2019. Credit card fraud classification based on GAN-AdaBoost-DT imbalanced classification algorithm. *Journal of Computer Applications* 39: 618–22.
23. Charitou, Charitos, Simo Dragicevic, and Artur d'Avila Garcez. 2021. Synthetic data generation for fraud detection using GANs. *arXiv arXiv:2109.12546*.
24. Scott, Mitchell, and Jo Plested. 2019. GAN-SMOTE: A generative adversarial network approach to synthetic minority oversampling. *Australian Journal of Intelligent Information Processing Systems* 15: 29–35.
25. Zhang, Liyuan, Huamin Yang, and Zhengang Jiang. 2018. Imbalanced biomedical data classification using self-adaptive multilayer ELM combined with dynamic GAN. *Biomedical Engineering Online* 17: 181.
26. Lan, Lan, Lei You, Zeyang Zhang, Zhiwei Fan, Weiling Zhao, Nianyin Zeng, Yidong Chen, and Xiaobo Zhou. 2020. Generative adversarial networks and its applications in biomedical informatics. *Frontiers in Public Health* 8: 164.
27. Bunkhumpornpat, Chumphol, Krung Sinapiromsaran, and Chidchanok Lursinsap. 2009. Safe-Level-SMOTE: Safe-level-synthetic minority over-sampling technique for handling the class imbalanced problem. In *Advances in Knowledge Discovery and Data Mining*. Edited by Thanaruk Theeramunkong, Boonserm Kijsirikul, Nick Cercone and Tu-Bao Ho. Berlin and Heidelberg: Springer, pp. 475–82.

28. Zhu, Xiaoqian, Xiang Ao, Zidi Qin, Yanpeng Chang, Yang Liu, Qing He, and Jianping Li. 2021. Intelligent financial fraud detection practices in post-pandemic era. *The Innovation* 2: 100176.
29. Ye, Huanzhuo, Lin Xiang, and Yanping Gan. 2019. Detecting financial statement fraud using random forest with SMOTE. *IOP Conference Series: Materials Science and Engineering* 612: 052051.
30. Jurgovsky, Johannes, Michael Granitzer, Konstantin Ziegler, Sylvie Calabretto, Pierre-Edouard Portier, Liyun He-Guelton, and Olivier Caelen. 2018. Sequence classification for credit-card fraud detection. *Expert Systems with Applications* 100: 234–45.
31. Branco, Bernardo, Pedro Abreu, Ana Sofia Gomes, Mariana S. C. Almeida, João Tiago Ascensão, and Pedro Bizarro. 2020. Interleaved sequence RNNs for fraud detection. Paper presented at the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20, New York, NY, USA, August 24–27; Melbourne: Association for Computing Machinery, pp. 3101–9.
32. Xie, Yu, Guanjun Liu, Chungang Yan, Changjun Jiang, and MengChu Zhou. 2022. Time-aware attention-based gated network for credit card fraud detection by extracting transactional behaviors. *IEEE Transactions on Computational Social Systems* 10: 1004–16.
33. Chen, Joy, and Kong-Long Lai. 2021. Deep convolution neural network model for credit-card fraud detection and alert. *Journal of Artificial Intelligence and Capsule Networks* 3: 101–12.
34. Zhang, Zhaojun, Xinxin Zhou, Xiaobo Zhang, Lizhi Wang, and Pengwei Wang. 2018. A model based on convolutional neural network for online transaction fraud detection. *Security and Communication Networks* 2018: 5680264.

Retrieved from <https://encyclopedia.pub/entry/history/show/117512>