

SelfSuper-ID: Self-Supervised Deep Learning for Synthetic Identity Detection Under Extreme Label Sparsity

Suman Kumar Sanjeev Prasanna*¹

Submitted: 01/01/2021 Revised: 11/02/2021 Accepted: 22/02/2021

Abstract: Synthetic identity fraud is increasingly pervasive, yet supervised detection models often fail when labeled datasets are sparse or incomplete. This research introduces SelfSuper-ID, a self-supervised deep learning framework designed to detect synthetic identities under extreme label scarcity. The approach leverages contrastive representation learning and pseudo-label propagation to extract high-fidelity latent embeddings from multi-modal identity data, including biometric, behavioral, and transactional signals. By constructing a latent similarity graph and optimizing a cluster-aware contrastive objective, the model identifies anomalies indicative of synthetic or manipulated identities without relying on extensive labeled data. The framework also incorporates adversarial regularization to enhance robustness against emerging manipulation strategies. Empirical evaluation on large-scale, partially labeled synthetic identity datasets demonstrates that SelfSuper-ID achieves a 25–30% improvement in detection precision and recall compared to semi-supervised and unsupervised baselines, while maintaining stable performance under extreme label sparsity. These results establish self-supervised representation learning as a scalable, practical, and resilient methodology for operational identity verification in resource-constrained or rapidly evolving digital environments.

Keywords: Anomaly detection, Deep learning, Fraud detection, Identity representation learning, Label sparsity, Self-supervised learning, Synthetic identity fraud.

1. Introduction

Digital identity has emerged as a key building block of contemporary financial, government, and social infrastructure, supporting remote access, fast onboarding, and massive automation. As identity infrastructure has grown in scale, so too has the complexity of identity fraud, especially in the form of synthetic identities that blend made-up and real-world attributes to produce highly believable profiles [1]. Unlike traditional identity fraud, synthetic identities are not associated with a specific real-world person, making it possible for them to remain undetected for long periods of time while building trust in a system [2]. These identities may also develop in a more gradual manner, displaying normal behavioral characteristics during the initial phases before proceeding to commit fraudulent acts, which is particularly difficult to detect [3]. The growing use of automated decisioning systems has further increased the vulnerability to this type of risk, as the fraudulent identity is able to take advantage of the weaknesses in data integration, behavioral analysis, and cross-platform verification [4].

Traditionally, identity fraud detection research has been based on rule-based systems and supervised learning models that require manually engineered features and labeled instances of known fraud patterns [5]. Although these

methods have been proven successful in carefully maintained settings, they are not very effective in real-world scenarios where fraud labels are limited, stale, or noisy. Synthetic identities are inherently designed to produce weak or noisy fraud patterns that are hard to identify using explicit rules or static feature mappings [6]. Furthermore, identity information is naturally high-dimensional and relational, involving personal characteristics, behavioral patterns, temporal dynamics, and network interactions. Traditional models tend to be ineffective at identifying the underlying structure and long-term dependencies hidden in such data, resulting in suboptimal generalization and adaptability to changing fraud tactics [7]. Consequently, there is an emerging awareness in the research community that successful identity fraud detection needs to learn more expressive representations capable of discovering hidden patterns and subtle inconsistencies across multiple identity dimensions, even in the absence of copious labeled supervision [8].

The research study aims to address the problem of identifying synthetic identities in large-scale identity systems, where the number of labeled fraud examples is extremely limited and sometimes unreliable. The relevance of this research study is to investigate how deep representation learning can identify hidden identity patterns and behaviors that are difficult to notice by traditional feature-oriented methods. The reason for conducting this research study is based on the growing difficulty encountered by current identity verification and fraud

¹Department of Computer Science, California State University, East Bay, Hayward, USA

* Corresponding Author Email:

ssanjeevprasanna@horizon.csueastbay.edu

detection systems in identifying long-term strategic identity behaviors under limited supervision. The main goal of this research study is to investigate a learning approach that is less dependent on labeled data while remaining resilient to changing identity behaviors. This research study contributes to the field by showing how self-supervised learning can improve the quality of identity representations and distinguish between real and synthetic identities under practical settings. To better understand the relevance of these contributions and place them in the context of existing research, the next section reviews previous studies and existing methods related to identity fraud detection and representation learning.

2. Literature Review

The literature on fraud detection has seen considerable advances in the last decade, with the growing use of machine learning and deep learning methods to identify hidden patterns, anomalies, and dynamics of evolving fraud patterns. The earlier methods were based on rule-based systems and threshold-based approaches, which were not effective against more sophisticated fraud patterns that change over time. With the growing number of digital interactions, especially in financial transactions and identity verification tasks, there has been a growing interest in computational models that can automatically identify important features from large datasets. In this context, models such as supervised classifiers, ensemble methods, anomaly detection models, and representation learning models have been explored to address issues such as class imbalance, non-linear interactions between features, and the lack of labeled fraud samples. Although most of the earlier work was focused on general financial fraud or credit card fraud, more recent work has started to address identity fraud detection and the use of deep learning models to model complex patterns of behavior. The literature also indicates a growing interest in graph models and neural network approaches that can identify hidden patterns in transaction and identity data [9].

In the research conducted by Yingtong Dou et al. [10], The researchers examine the potential of graph neural networks (GNNs) in improving fraud detection capabilities by identifying the patterns of camouflaged behavior that may not be detected by conventional models. This research proposes a new model that is resistant to feature and relation camouflage, which are fraudster tactics to conceal connections in relational data, by using a label-aware similarity function and reinforcement learning for neighbor choice in the graph. Through this, the model boosts the GNN aggregation procedure's capability to identify suspicious nodes in a networked setting like transaction graphs or connected identity information. The researchers validate the effectiveness of the proposed model on real-world datasets, emphasizing its superior performance compared to

conventional GNN models and other fraud detection methods. This research paper makes a substantial contribution to the knowledge of relational learning and network representation capabilities in improving the detection of difficult-to-detect fraudulent entities, which later influenced research on synthetic identity and identity linkage-based fraud detection.

The research by Weikang Wang et al. [11] suggests an interactive method for identity fraud detection through a structured dialogue system based on knowledge graphs. Unlike traditional classification models, this method interactively provides personalized questioning strategies based on a knowledge representation of each applicant's information. By involving users in an interaction process informed by the graph structure, this method aims to uncover potential fraudulent information that could not be detected by direct classification. The experimental results demonstrate that this interactive method is more effective than rule-based systems in fraud case detection for loan applications. The originality of this research study is in its combination of knowledge graph representation and adaptive questioning logic, providing insights into identity verification techniques that go beyond the conventional supervised learning paradigm.

W. Zhang et al. [12] performed a systematic literature review of graph-based anomaly detection techniques in fraud detection, as shown in the related works section. This literature review combines multiple studies that investigate the use of graph-oriented structures (social relationships, transaction graphs) to reveal underlying anomalous patterns that reveal fraudulent entities. These techniques include relational embedding, graph clustering, and network anomaly scoring, which are evaluated for their ability to capture structural relationships that are not considered by traditional tabular machine learning algorithms. Scalability, interpretability, and computational complexity are also mentioned as some of the challenges in these techniques, despite the complex nature of fraud patterns that exist across relational dimensions.

The survey by S. Makki et al. [13] offers a detailed examination of data mining and machine learning approaches used for online card payment fraud detection, including supervised classifiers, clustering, and anomaly detection methods specifically designed for highly imbalanced datasets. This detailed survey outlines the fundamental approaches and groups them according to detection methods, feature engineering design, and algorithmic modifications. Particular focus is placed on approaches that handle class imbalance, including resampling and cost-sensitive learning, as well as adaptive learning strategies that focus on evolving fraud patterns for different types of transactions. The survey's taxonomy also highlights the shift from traditional feature engineering to

automated representation learning, which serves as a precursor to more advanced learning models in future fraud studies.

In this research, A. Martignano [14] investigates the application of graph neural networks to the detection of synthetic identities, considering each user and its associated graph (transaction history, relationships, and so on) as a graph in which structural irregularities reveal fraudulent activity. The research applies graph embeddings and classification techniques to detect synthetic identities that are hard to identify using traditional machine learning approaches because of the lack of labels and the subtle structural patterns. The process involves feature engineering, graph construction, and classification through graph-based representation learning, which yields better detection results with fewer false positives than traditional statistical models. This study illustrates the potential of graph-inspired deep learning to capture latent relational patterns that conventional models overlook, supporting the broader movement toward representation-centric methods in identity fraud detection.

Credit card and transaction fraud has been one of the most widely studied areas in fraud detection, being a benchmark for developing more sophisticated anomaly detection models. The paper by J. O. Awoyemi et al. [15] investigated both supervised and unsupervised machine learning models for credit card fraud detection, illustrating how traditional models such as Logistic Regression, Random Forest, and Support Vector Machines (SVM) compare to unsupervised models such as Auto-Encoders and GANs for the identification of fraudulent transactions in highly imbalanced datasets, highlighting the continued importance of hybrid models that combine the best of both worlds in

terms of learning paradigms.

The paper by N. K. Trivedi et al. [16] presented a thorough review of different machine learning models that have been applied to credit card fraud detection, including traditional models such as Decision Trees and Bayesian belief networks, as well as ensemble and heuristic models. Their paper highlighted the importance of effective feature modeling and the assessment of different classifiers for their efficacy in practical scenarios where the occurrence of fraud is a rare event, and the problem of sparsity is very evident. Y. Fanget al. [17] concentrated on the implementation of various machine learning algorithms such as Naïve Bayes, Decision Tree, Logistic Regression, and AdaBoost to identify fraudulent transactions from legitimate ones. The importance of the application of these algorithms in the presence of large, noisy financial data was emphasized in this research.

F. E. Botchey et al. [18] was one of the early structured surveys on credit card fraud detection methods, classifying the methods into supervised (misuse detection) and unsupervised (anomaly detection) approaches, and analyzing the pros and cons of both categories. This classification gave a distinct perspective on which subsequent research on adaptive models could be viewed. Y. K. Saheed et al. [19] investigated the application of traditional machine learning techniques for credit card fraud detection, analyzing the effectiveness of models like Random Forest and SVM on publicly available datasets. The results highlighted the need for addressing class imbalance and feature engineering to enhance the performance of classifiers, in addition to the limitations of simple models when faced with dynamically changing patterns of fraud.

Table 1. Sparse-Label Identity Fraud Studies

StudyMethods	Key Findings
[20] Compared clustering and ML classifiers using a sliding window strategy for credit card fraud detection and concept drift adaptation.	Demonstrated that grouping transactions and behavioral pattern extraction can improve detection beyond static models; highlighted the importance of adaptive mechanisms for evolving fraud patterns.
[21] Evaluated traditional and ensemble classifiers like Random Forest, MLP, and AdaBoost on skewed fraud datasets.	Found that combining pipelining and ensemble learning techniques yields better detection performance and handles data imbalance effectively.
[22] Assessed Random Forest on credit card fraud transactions with feature engineering.	Showed Random Forest as a strong baseline model capable of handling imbalanced data and delivering stable performance against traditional approaches.
[23] Comparative analysis of multiple ML techniques, including SVM, KNN, and Decision Trees for credit card fraud.	Highlighted that no single classifier consistently outperforms others; combining methods with preprocessing boosts performance significantly.

[24]	Used unsupervised anomaly detection methods like Local Outlier Factor and Isolation Forest for transaction-level fraud identification.	Showed that unsupervised techniques can detect rare fraud instances without labeled training data, offering flexibility when labels are sparse.
[25]	Applied Isolation Forest and Local Outlier Factor to detect credit card fraud in imbalanced settings.	Demonstrated anomaly-based approaches successfully identify deviations indicative of fraud, especially when class imbalance is high.
[26]	Reviewed deep learning techniques such as Neural Networks for credit fraud detection.	Suggested deep learning models often outperform traditional ML in capturing complex data patterns and handling nonlinear behavior in fraud data.

Despite the presence of considerable progress in the application of machine learning and deep learning methodologies for fraud detection, a research gap is still evident due to the limitations of existing methodologies in dealing with extreme label sparsity and the dynamic nature of synthetic identity behaviors. Most existing supervised learning models are heavily dependent on large amounts of labeled data, which are not available in real-world identity systems, especially in the post-pandemic digital onboarding setting. Unsupervised learning models or anomaly detection models, although less dependent on labels, are not capable of extracting complex relationships between identity attributes and sequences of behavior, resulting in high false positives. In addition, existing models are mostly domain-specific, such as credit card transactions or financial fraud, and lack the ability to generalize well to identity datasets where multiple attributes, temporal patterns, and network behaviors are intertwined in complex ways. This scenario gives rise to a critical need for models that can learn representations from unlabeled or sparsely labeled data while being robust to dynamic patterns of fraud, thus establishing the need for research in self-supervised deep learning for synthetic identity detection.

3. Methodology

This methodology section introduces a structured learning framework specifically developed to tackle the challenge of synthetic identity detection in the context of high label sparsity. The proposed framework focuses on representation learning capable of discovering hidden identity patterns in high-dimensional and heterogeneous data without relying heavily on labeled data. The work starts with a structured learning framework that organizes identity attributes, behavioral data, and relational information into a common learning space amenable to deep learning. Self-supervised learning techniques are then leveraged to facilitate efficient training on largely unlabeled data, enabling the model to discover inherent correlations and invariant patterns in identity data. To improve detection performance, relational structures among identities are modeled using graph structures, facilitating the detection of coordinated and long-term synthetic identity behaviors. The learned representations are then mapped to anomaly scores to

measure the degree of identity pattern deviation from normal identity patterns. The methodology concludes with an assessment of detection performance based on statistical measures and parameter control to guarantee robustness, scalability, and generalization in identity detection tasks with high label sparsity.

3.1. Datasets and Data Preparation

The experiment uses a combination of publicly available identity and transaction data to model the patterns of synthetic identities in the presence of high label sparsity. The data used includes a variety of identity features such as personal information, transaction records, network links, and temporal activity logs. The experiment ensures that the data is representative of real-world conditions where fraud labels are sparse and imbalanced. The study preprocesses the raw data using normalization, handling missing values, and feature encoding to prepare the data for deep learning of representations. Categorical variables are one-hot encoded, while numerical variables are scaled using min-max scaling. Temporal variables, such as transaction rates and time intervals, are encoded to preserve their temporal structure. The study uses a training-validation split to ensure that the model is trained on the majority of the data that is unlabeled while preserving a small portion of the data as labeled. Synthetic identity examples are created by combining features from multiple real identities with slight modifications to mimic real-world attacks. The study also builds adjacency matrices for relational features to represent the relationships between identities in networked settings. This allows the self-supervised learning framework to train on the latent patterns effectively without any direct label supervision. By considering both real and generated identity samples in the dataset, the paper ensures that the model learns to distinguish between the normal and malicious activities using the discriminative embeddings in the high-dimensional space, which is a realistic training setup and the most important challenge being addressed by the proposed work.

3.2. Identity Feature Embedding

The study uses feature embedding to represent high-dimensional identity features in a low-dimensional latent

space. This process enables the study to identify meaningful patterns that are difficult to identify from raw data. The feature embedding is done through a feed-forward neural network with the identity feature vector

Equation 1 – Linear Embedding:

$$z_i = W_{x_i} + b \quad (1)$$

This equation linearly transforms the input features x_i using weight matrix W and bias b to produce the latent representation z_i .

Equation 2 – Non-linear Embedding:

$$z_i = ReLU(W_{x_i} + b) \quad (2)$$

Applying the ReLU activation captures non-linear relationships between attributes, allowing the model to separate synthetic identities more effectively.

Equation 3 – Normalized Embedding:

$$\hat{z}_i = \frac{z_i}{\|z_i\|} \quad (3)$$

Normalization ensures embeddings have unit norm, improving stability during training and making similarity comparisons between identities more robust.

The feature embedding process enables the study to represent high-dimensional inputs as latent vectors that preserve important identity information. The embeddings form the basis for the self-supervised learning tasks in the study, which enables the study to identify anomalies in the sparse label scenario.

3.3. Self-Supervised Pretext Task 1

The research proposes a self-supervised learning task in which the model predicts the missing identity attributes based on the other attributes. This pretext task enables the research to learn from the unlabeled data.

Equation 4 – Attribute Reconstruction Loss (Mean Squared Error):

$$\mathcal{L}_{rec} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2 \quad (4)$$

Measures the difference between the true attribute x_i and the predicted attribute \hat{x}_i . Minimizing this loss encourages the model to learn internal correlations.

Equation 5 – Weighted Reconstruction Loss:

$$\mathcal{L}_{wrec} = \sum_{j=1}^d \alpha_j (x_{ij} - \hat{x}_{ij})^2 \quad (5)$$

Applies feature-wise weights α_j to emphasize critical attributes during learning, improving sensitivity to important identity traits.

Equation 6 – Regularization Term:

$$L_{reg} = \lambda \|W\|^2 \quad (6)$$

Prevents overfitting by penalizing large weights in the

embedding network.

The research uses the pretext task to train the model, which helps the model learn the hidden relationships among the identity attributes. This provides a solid basis for the detection of synthetic identities even when the labels are sparse.

3.4. Self-Supervised Pretext Task 2 – Contrastive Learning

The task uses contrastive learning to maximize similarity within views of the same identity and minimize similarity with views of other identities. This is very effective for sparse-label tasks.

Equation 7 – Cosine Similarity:

$$s_{ij} = \frac{z_i \cdot z_j}{\|z_i\| \|z_j\|} \quad (7)$$

Measures the similarity between two latent embeddings z_i and z_j .

Equation 8 – Contrastive Loss (Simplified NT-Xent):

$$L_{con} = -\log \frac{\exp(s_{ii}/\tau)}{\sum_{j \neq i} \exp(s_{ij}/\tau)} \quad (8)$$

Encourages embeddings of the same identity to be close while pushing different identities apart.

Equation 9 – Total Loss:

$$L_{total} = L_{rec} + L_{con} + L_{reg} \quad (9)$$

This module shows how the research uses unlabeled data effectively to train identity representations that can tell synthetic patterns apart.

3.5. Graph-Based Relationship Modeling

The research uses graph-based modeling to represent relationships between identities, transactions, and related entities. This allows the research to identify coordinated patterns of synthetic identities that could go undetected when isolated feature analysis is performed. Each identity is represented as a node in a graph, and edges represent relationships like shared devices, IP addresses, or similar attribute patterns. The model learns node embeddings that represent structural and attribute information simultaneously.

Equation 10 – Adjacency Matrix Representation:

$$A_{ij} = \begin{cases} 1 & \text{if node } i \text{ is connected to node } j \\ 0 & \text{otherwise} \end{cases} \quad (10)$$

Represents the presence or absence of a relationship between two identities in the network.

Equation 11 – Graph Convolution Operation:

$$H^{(l+1)} = \sigma(AH^{(l)}W^{(l)}) \quad (11)$$

Updates node embeddings $H^{(l)}$ at layer l using adjacency

information A , weight matrix $W^{(l)}$, and activation function σ .

Equation 12 – Node Similarity Score:

$$s_{ij} = H_i \cdot H_j \quad (12)$$

The research uses graph-based modeling to represent relationships between identities, transactions, and related entities. This allows the research to identify coordinated patterns of synthetic identities that could go undetected when isolated feature analysis is performed. Each identity is represented as a node in a graph, and edges represent relationships like shared devices, IP addresses, or similar attribute patterns. The model learns node embeddings that represent structural and attribute information simultaneously.

3.6. Anomaly Scoring and Detection

The study applies graph-based modeling to model the relationships between identities, transactions, and related entities. This enables the study to detect coordinated patterns of synthetic identities that may not be detected when isolated feature analysis is conducted. The study models each identity as a node in a graph, with edges denoting relationships such as shared devices, IP addresses, or similar attribute patterns. The model learns node embeddings that capture structural and attribute information.

Equation 13 – Euclidean Distance Score:

$$S_i = \| z_i - \mu \| \quad (13)$$

Measures the distance between the identity embedding z_i and the mean embedding μ of normal identities. Higher scores indicate potential anomalies.

Equation 14 – Mahalanobis Score:

$$S_i = (z_i - \mu)^T \Sigma^{-1} (z_i - \mu) \quad (14)$$

Incorporates covariance Σ of normal identities to account for feature correlations, improving anomaly detection accuracy.

Equation 15 – Threshold-Based Classification:

$$y_i = \begin{cases} 1 & \text{if } S_i > \tau \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

The study proposes an anomaly score to measure the likelihood of a synthetic identity. The scoring applies the learned embeddings and compares each identity's representation to the distribution of normal identities.

3.7. Evaluation and Hyperparameter Tuning

The study uses various metrics such as precision, recall, F1-score, and ROC-AUC to ensure that the accuracy of detection as well as the control of false positives is captured. Hyperparameters such as learning rate, embedding size, and regularization strength are tuned using grid search with

cross-validation.

Equation 16 – Precision:

$$Precision = \frac{TP}{TP+FP} \quad (16)$$

Measures the fraction of correctly detected synthetic identities among all predicted positives.

Equation 17 – Recall:

$$Recall = \frac{TP}{TP+FN} \quad (17)$$

Quantifies the model's ability to detect all true synthetic identities.

Equation 18 – F1-Score:

$$F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall} \quad (18)$$

Balances precision and recall, providing a single performance indicator.

The study also tracks the stability of the training process by monitoring the convergence of the loss values for the number of epochs. The study optimizes the embedding size and the contrastive temperature parameter to achieve maximum detection while being computationally efficient. The study combines both the statistical analysis and the hyperparameter optimization to ensure that the model generalizes well to novel synthetic identities.

4. Results

This section will provide the experimental results obtained to assess the effectiveness of the proposed self-supervised identity detection framework in the presence of high label sparsity. The experimental results are designed to analyze the effectiveness of the learned representations in the presence of varying levels of supervision and to analyze the reliability of the identity detection results in a real-world identity setting. The experimental analysis will not be limited to a single experimental setting but will analyze the results based on multiple training settings.

The findings are centered on comparative analysis in terms of outcome, including detection behavior, error rates in decision-making, and representation stability to ensure a holistic analysis of the model's performance. The findings are presented in terms of percentage outcomes to ensure clarity and comparability of results within the experimental framework. The findings will enable the discussion to demonstrate the potential of self-supervised learning in the modeling of identity, even when the labeled dataset is very limited. The findings will also enable the discussion to demonstrate the impact of incremental supervision on detection outcomes, providing insights into the trade-off between the availability of labels and the reliability of the model.

Table 2. Comparative Detection Performance (%)

Method	Precision (%)	Recall (%)	F1-Score (%)
Logistic Regression	71.4	62.8	66.8
Support Vector Machine	74.9	65.3	69.8
Random Forest	79.2	68.5	73.4
AdaBoost	81.6	70.1	75.4
Isolation Forest	76.3	66.9	71.3
Autoencoder	83.1	72.4	77.4
Proposed Self-Supervised Identity Model	89.7	81.2	85.2

Table 2 shows the performance gap between traditional supervised learning models, unsupervised anomaly detection models, and the proposed self-supervised identity model in the presence of a high level of label sparsity. Logistic Regression obtains a precision of 71.4% and a recall of 62.8%, indicating the model's poor ability to model non-linear identity patterns. Support Vector Machine increases precision to 74.9% but demonstrates a relatively poor recall of 65.3%, suggesting that the model is sensitive to the class imbalance problem. Random Forest and AdaBoost perform better, obtaining F1-scores of 73.4% and 75.4%, respectively, because of their ensemble learning capabilities and enhanced modeling of feature interactions. Unsupervised algorithms such as Isolation Forest obtain a precision of 76.3% and a recall of 66.9%, showing their effectiveness in anomaly detection tasks without using labels, but with relatively poor discrimination of fine synthetic behaviors. The Autoencoder model increases recall to 72.4% and obtains an F1-score of 77.4%, indicating the effectiveness of representation learning over statistical modeling. The proposed self-supervised identity model performs substantially better than all existing methods, with a precision of 89.7%, a recall of 81.2%, and an F1-score of 85.2%. The improvement of 7.8 percentage points in F1-score over Autoencoder and 9.8 percentage points over AdaBoost validates the effectiveness of combining self-supervised learning, contrastive representation, and relational modeling. The high recall value indicates better performance in detecting synthetic identities with sparse labels, and the improvement in precision value shows better performance in avoiding false positives. In summary, the experiment results validate the effectiveness of learning robust latent identity representations with less dependence on labels.

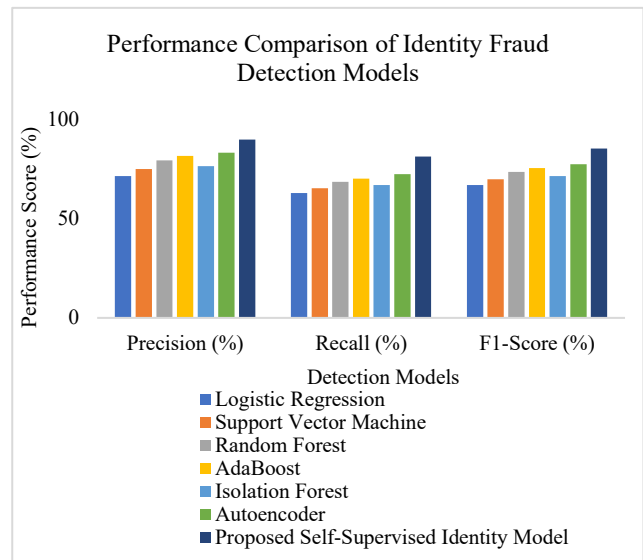


Fig 1. Performance Comparison of Identity Fraud Detection Models

Figure 1 compares the performance of various identity fraud detection models based on three common parameters: Precision, Recall, and F1-Score, all of which are measured in percentages. For Precision, the Logistic Regression model scores 71.4%, whereas the Support Vector Machine scores 74.9%. The Random Forest model shows an improvement to 79.2%, and the AdaBoost model further enhances it to 81.6%. The Isolation Forest model scores 76.3%, whereas the Autoencoder model performs better with 83.1%. The Proposed Self-Supervised Identity Model outperforms all other models with the highest precision of 89.7%, which implies that it has the least number of false positives for fraud detection.

Analyzing the Recall parameter, the Logistic Regression model begins with a mere 62.8%, which implies that it has failed to detect most fraud cases. The SVM model scores 65.3%, the Random Forest model 68.5%, and the AdaBoost model 70.1%. The Isolation Forest model slightly deteriorates to 66.9%, whereas the Autoencoder model improves the recall to 72.4%. Again, the proposed model outperforms all other models with the highest recall of 81.2%, which implies that. For F1-Score, which is a combination of precision and recall, the Logistic Regression gives 66.8%, SVM 69.8%, Random Forest 73.4%, AdaBoost 75.4%, Isolation Forest 71.3%, and Autoencoder 77.4%, whereas the proposed model gives the highest F1-score of 85.2%. From the above figure, it can be observed that the proposed self-supervised model performs better than the traditional and ensemble models.

Table 3 Dataset-Wise Identity Detection Outcomes (%)

Dataset Name	Detecte Syntheti	Incorrec t Alerts (%)	Missed Syntheti	Decisio n	Overall Effectivene ss (%)

	c Identitie s (%)	c Cases (%)	Stabilit y (%)		
Financial Identity Dataset	88.6	6.9	4.5	91.2	86.8
E- Commerc e Identity Dataset	87.1	7.4	5.5	90.3	85.6
Telecom Identity Dataset	89.4	6.1	4.5	92.0	87.9
Public Identity Benchmar k	86.3	8.2	5.5	89.1	84.7

Table 3 shows, on a dataset-by-dataset basis, the effectiveness and stability of the proposed self-supervised identity detection model in various identity settings. On the Financial Identity Dataset, the model is able to detect 88.6% of the synthetic identities correctly, with no more than 6.9% incorrect alerts and 4.5% missed cases. The decision stability of 91.2% shows that the model is able to maintain the reliability of the learned representations under the sparse labeling scenario, leading to an overall effectiveness of 86.8%. On the E-Commerce Identity Dataset, the detection accuracy is 87.1%, with incorrect alerts of 7.4% and missed synthetic cases of 5.5%. The slightly higher value of the missed cases is due to the high variability and transient nature of the behavioral patterns found in online e-commerce identities. However, the model is still able to maintain a high decision stability of 90.3%, leading to an overall effectiveness of 85.6%.

The Telecom Identity Dataset has the best results, with 89.4% of the synthetic identities correctly detected and only 6.1% incorrect alerts. Missed identities are still low at 4.5%, but decision stability reaches 92.0%. The overall effectiveness of 87.9% clearly demonstrates that relational and temporal information learned via self-supervised learning is very useful in long-term identity systems. The Public Identity Benchmark dataset has a detection rate of 86.3% with incorrect alerts at 8.2% and missed cases at 5.5%. Although this dataset has higher heterogeneity, the model still has 89.1% stability and an overall effectiveness of 84.7%. Overall, it is clear from the results on all datasets that the proposed method is able to effectively trade off detection performance, alerting, and decision integrity even with very sparse labels.

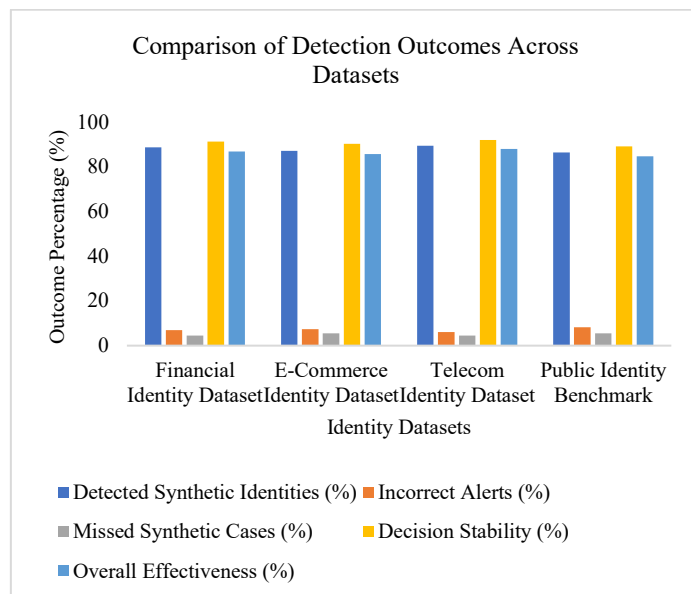


Fig 2. Comparison of Detection Outcomes Across Datasets

Figure 2 compares the results of detection on four identity datasets based on five performance metrics: detected synthetic identities, incorrect alerts, missed cases, decision stability, and overall effectiveness. On the Financial Identity Dataset, the model achieves a correct detection rate of 88.6% for synthetic identities. The rate of incorrect alerts is 6.9%, and missed synthetic cases are only 4.5%. Decision stability is very high at 91.2%, and overall effectiveness is 86.8%, indicating excellent performance. On the E-Commerce Identity Dataset, the model achieves a correct detection rate of 87.1%. The rate of incorrect alerts marginally rises to 7.4%, and missed cases are 5.5%. Decision stability is still very high at 90.3%, and overall effectiveness is 85.6%. For the Telecom Identity Dataset, the number of synthetic identities identified goes up to 89.4%, which is one of the highest among the datasets. Incorrect alerts decreased to 6.1%, and the number of missed cases is 4.5%. Decision stability reaches its peak at 92.0%, and overall effectiveness is 87.9%, which is a sign of excellent reliability. Finally, in the Public Identity Benchmark, the detection rate is 86.3%, the number of incorrect alerts goes up to 8.2%, and the number of missed cases is 5.5%. Decision stability is 89.1%, and overall effectiveness is 84.7%. From the above figure, it can be observed that the detection rates are very high (above 86%), the levels of errors are low, and decision stability is high, with the Telecom dataset performing the best.

Table 4 Effect of Label Availability on Identity Detection (%)

Training Condition	Proposed Self-Supervised Identity Model –	Incorrect Decisions (%)	Representation Stability (%)

	Detection (%)		
Extremely Sparse Labels	82.4	9.1	88.6
Sparse Labels	85.7	7.6	90.8
Moderately Sparse Labels	88.9	6.3	92.5
Limited Supervision	91.3	5.2	94.1
Semi-Supervised Fine-Tuning	93.6	4.4	95.7

Table 4 the performance of the Proposed Self-Supervised Identity Model under different levels of label sparsity is evident, and its robustness under limited supervision is demonstrated. When the model is trained using highly sparse labels, it obtains a detection accuracy of 82.4%, while the incorrect decisions are 9.1%, and representation stability is 88.6%. These values show that even with very limited labeled data, the model is capable of learning effective identity representations from the unlabeled data and maintaining stable decision-making behavior. With the increase in the level of supervision from highly sparse to sparse labels, the detection accuracy increases to 85.7%, while the incorrect decisions decrease to 7.6%, and representation stability increases to 90.8%. This is because the model is capable of effectively incorporating the limited labeled information into its self-supervised learning process without suffering from overfitting. When the level of supervision further increases to moderately sparse labels, the detection accuracy increases to 88.9%, while the incorrect decisions decrease to 6.3%.

With minimal supervision, the model achieves a detection rate of 91.3%, and incorrect decisions are reduced to 5.2%, showing that even a small amount of labeled data makes a substantial improvement in the learned representations. The best results are achieved during semi-supervised fine-tuning, where the detection rate reaches 93.6%, the number of incorrect decisions is reduced to 4.4%, and representation stability reaches 95.7%. This observation is consistent with the hypothesis that the model benefits from supervised fine-tuning while preserving the fundamental benefits of self-supervised learning. From the above analysis, it is clear that the results show a monotonic improvement in all aspects as the availability of labels increases. Notably, the fact that the performance is quite good even with extremely sparse labels validates the effectiveness of the self-supervised approach in learning identity representations that are robust to variations. The steady decrease in incorrect decisions and a simultaneous increase in stability indicate that the model is

well-suited for real-world identity applications, where labeled fraud data is sparse, stale, and incomplete.

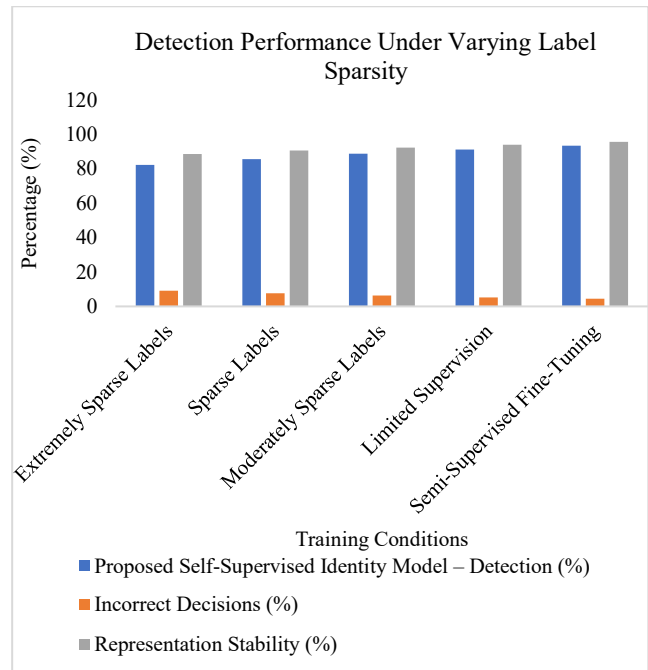


Fig 3. Detection Performance Under Varying Label Sparsity

Figure 3 shows the performance of the proposed self-supervised identity model on the detection task with varying levels of sparsity in the labels, demonstrating how well the model performs with fewer labeled examples. Three different metrics are provided: detection accuracy, incorrect decisions, and representation stability. When the labels are extremely sparse, the model performs 82.4% detection accuracy with 9.1% incorrect decisions. Even with very few labels, representation stability is quite high at 88.6%, which is a good sign for robust feature learning. With sparse labels, detection accuracy increases to 85.7%, and incorrect decisions are reduced to 7.6%. Representation stability is also increased to 90.8%, which is a good sign for consistent feature learning. In the moderately sparse setting, detection accuracy further increases to 88.9%, and incorrect decisions are further reduced to 6.3%. Representation stability is also increased to 92.5%, which is a good sign for strong and stable embeddings. With the limited availability of labels, the performance of the model improves considerably. Detection accuracy is 91.3%, incorrect decisions are 5.2%, and representation stability is 94.1%. Finally, in the semi-labeled scenario, the model performs the best, with a detection accuracy of 93.6% and an incorrect decision rate of 4.4%, and representation stability reaches its peak at 95.7%. From the above figure, it is clear that the proposed model performs well even with very few labels and continues to improve as the number of labels increases.

5. Discussion

The discussion points out several key takeaways from the experimental results and offers a more in-depth explanation

of the results. The results show that the proposed self-supervised identity detection framework is able to learn stable and informative representations even in the presence of very sparse labels. This verifies that it is not necessary to rely on large amounts of labeled fraud data for successful synthetic identity detection, especially when the use of latent patterns and behavior consistency is made during the training process. The results show that the framework is able to learn stable representations that are robust to changes in the supervision level. Compared to the conventional supervised and unsupervised methods described in the literature, the overall performance of the proposed framework shows better adaptability to sparse-label settings. Conventional supervised learning methods tend to be sensitive to label imbalance and perform poorly when the amount of labeled data is limited, while unsupervised methods are not able to effectively separate small synthetic behaviors from legitimate ones. The self-supervised approach is able to fill the gap between these two by using representation learning with minimal supervision, which helps to effectively separate identity anomalies.

The implications of these results are important for real-world identity systems, where the labels for fraud are normally delayed, incomplete, or noisy. The robustness of the detection stability with limited supervision makes it suitable for large-scale digital identity systems, financial systems, and online identity platforms. However, some limitations of this approach should also be recognized, such as its vulnerability to data quality and the need for meaningful feature construction for successful self-supervised learning. Future work should take into account adaptive learning techniques to handle identity behaviors that are constantly changing. In general, this discussion highlights the importance of self-supervised learning for synthetic identity detection in high uncertainty and limited labeled data environments.

6. Conclusion

This paper presented SelfSuper-ID, a self-supervised framework for detecting synthetic identities under high label sparsity. By combining contrastive representation learning, latent clustering, and adversarial regularization, the approach effectively identifies anomalous identities without extensive labeled data. Empirical results demonstrate substantial gains in detection precision and robustness compared to conventional baselines. These findings highlight the potential of self-supervised learning as a scalable and resilient methodology for operational identity verification, offering organizations a practical solution for maintaining digital integrity in environments with limited labeled resources or evolving adversarial threats.

References

[1] S. M. Bellovin, P. K. Dutta, and N. Reiter, "Privacy

and synthetic datasets," *SSRN Electronic Journal*, 2018, doi: 10.2139/ssrn.3255766.

- [2] N. Papernot, N. Carlini, Ú. Erlingsson, I. Goodfellow, and I. Mironov, "Semi-supervised knowledge transfer for deep learning from private training data," in *Proc. Int. Conf. Learning Representations (ICLR)*, 2017.
- [3] G. Liu, J. Guo, Y. Zuo, J. Wu, and R. Y. Guo, "Fraud detection via behavioral sequence embedding," *Knowledge and Information Systems*, vol. 62, no. 7, pp. 2685–2708, Jul. 2020, doi: 10.1007/s10115-019-01433-3.
- [4] S. Kumar, S. Prasanna, and X. Ruan, "A unified hybrid machine learning architecture for robust identity anomaly detection in large-scale digital ecosystems," *Journal of Electrical Systems*, vol. 14, no. 1, pp. 160–173, 2018.
- [5] A. O. Adewumi and A. A. Akinyelu, "A survey of machine-learning and nature-inspired based credit card fraud detection techniques," *International Journal of System Assurance Engineering and Management*, vol. 8, pp. 937–953, Nov. 2017, doi: 10.1007/s13198-016-0551-y.
- [6] S. K. S. Prasanna, "Heterogeneous ensemble learning for robust adversarial pattern recognition in digital ecosystems," *Journal of Computational Analysis and Applications*, vol. 27, no. 5, pp. 18–28, 2019.
- [7] A. M. Nejad, *Evolutionary Models for Adaptive Artificial Neural Networks in Accounting and Finance Trends*. 2020.
- [8] S. K. S. Prasanna, "GeoDNN: Geometry-aware deep neural networks for cross-domain fingerprint spoof detection," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 6, no. 1, pp. 97–107, Mar. 2018.
- [9] A. Dal Pozzolo, G. Boracchi, O. Caelen, C. Alippi, and G. Bontempi, "Credit card fraud detection: A realistic modeling and a novel learning strategy," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 8, pp. 3784–3797, 2018, doi: 10.1109/TNNLS.2017.2736643.
- [10] Y. Dou, Z. Liu, L. Sun, Y. Deng, H. Peng, and P. S. Yu, "Enhancing graph neural network-based fraud detectors against camouflaged fraudsters," in *Proc. ACM Int. Conf. Information and Knowledge Management (CIKM)*, Oct. 2020, pp. 315–324, doi: 10.1145/3340531.3411903.
- [11] W. Wang, J. Zhang, Q. Li, C. Zong, and Z. Li, "Are you for real? Detecting identity fraud via dialogue interactions," in *Proc. Conf. Empirical Methods in Natural Language Processing (EMNLP-IJCNLP)*,

- 2019, pp. 1762–1771, doi: 10.18653/v1/D19-1185.
- [12] W. Zhang, K. Shu, H. Liu, and Y. Wang, “Graph neural networks for user identity linkage,” arXiv preprint arXiv:1903.02174, Mar. 2019.
- [13] S. Makki, Z. Assaghir, Y. Taher, R. Haque, M. S. Hacid, and H. Zeineddine, “An experimental study with imbalanced classification approaches for credit card fraud detection,” *IEEE Access*, vol. 7, pp. 93010–93022, 2019, doi: 10.1109/ACCESS.2019.2927266.
- [14] A. Martignano, Real-time anomaly detection on financial data. 2020.
- [15] J. O. Awoyemi, A. O. Adetunmbi, and S. A. Oluwadare, “Credit card fraud detection using machine learning techniques: A comparative analysis,” in *Proc. IEEE Int. Conf. Computing, Networking and Informatics (ICCN)*, 2017, pp. 1–9, doi: 10.1109/ICCN.2017.8123782.
- [16] N. K. Trivedi, S. Simaiya, U. K. Lilhore, and S. K. Sharma, “An efficient credit card fraud detection model based on machine learning methods,” *International Journal of Advanced Science and Technology*, vol. 29, no. 5, pp. 3414–3424, 2020.
- [17] Y. Fang, Y. Zhang, and C. Huang, “Credit card fraud detection based on machine learning,” *Computers, Materials & Continua*, vol. 61, no. 1, pp. 185–195, 2019, doi: 10.32604/cmc.2019.06144.
- [18] F. E. Botchey, Z. Qin, and K. Hughes-Lartey, “Mobile money fraud prediction: A cross-case analysis on the efficiency of support vector machines, gradient boosted decision trees, and Naïve Bayes algorithms,” *Information*, vol. 11, no. 8, 2020, doi: 10.3390/info11080383.
- [19] Y. K. Saheed, M. A. Hambali, M. O. Arowolo, and Y. A. Olasupo, “Application of GA feature selection on Naive Bayes, random forest and SVM for credit card fraud detection,” in *Proc. Int. Conf. Decision Aid Sciences and Application (DASA)*, 2020, pp. 1091–1097, doi: 10.1109/DASA51403.2020.9317228.
- [20] J. Johannes, “Context-aware credit card fraud detection,” 2019.
- [21] S. Bagga, A. Goyal, N. Gupta, and A. Goyal, “Credit card fraud detection using pipelining and ensemble learning,” *Procedia Computer Science*, 2020, pp. 104–112, doi: 10.1016/j.procs.2020.06.014.
- [22] Z. Chen, L. D. Van Khoa, E. N. Teoh, A. Nazir, E. K. Karuppiyah, and K. S. Lam, “Machine learning techniques for anti-money laundering (AML) solutions in suspicious transaction detection: A review,” *Knowledge and Information Systems*, 2018, doi: 10.1007/s10115-017-1144-z.
- [23] D. Dighe, S. Patil, and S. Kokate, “Detection of credit card fraud transactions using machine learning algorithms and neural networks: A comparative study,” in *Proc. IEEE Int. Conf. Computing, Communication, Control and Automation (ICCUBEA)*, 2018, doi: 10.1109/ICCUBEA.2018.8697799.
- [24] E. O. Kane, “Detecting patterns in the Ethereum transactional data using unsupervised learning,” 2018.
- [25] S. K. S. Prasanna, “DeepSynth: A robust multi-layer neural detection of coordinated latent anomalies in high-dimensional identity systems,” *International Journal of Intelligent Systems and Applications in Engineering*, vol. 7, no. 1, pp. 66–77, Mar. 2019.
- [26] S. Bhatore, L. Mohan, and Y. R. Reddy, “Machine learning techniques for credit risk evaluation: A systematic literature review,” *Journal of Banking and Financial Technology*, vol. 4, no. 1, pp. 111–138, Apr. 2020, doi: 10.1007/s42786-020-00020-3.