

Data Mining Project Report
Ali Raza - 20I-0782 | Hammad Umar - 19I-2157
Section A

Table of Contents

| | |
|--|----------|
| Table of Contents..... | 2 |
| Introduction:..... | 3 |
| Data Exploration and Preprocessing:..... | 3 |
| Data Augmentation with Geometric Distribution Mask..... | 3 |
| Model Architecture:..... | 3 |
| Model Training:..... | 4 |
| Evaluation Metrics:..... | 4 |
| Comparison with Traditional Methods:..... | 4 |
| Results and Discussion:..... | 4 |
| Conclusion:..... | 4 |

Introduction:

The detection of fraudulent transactions in credit card data is crucial for financial institutions to prevent losses and maintain customer trust. In this report, we present an approach for anomaly detection in credit card transactions using various techniques including traditional methods and deep learning-based models.

Data Exploration and Preprocessing:

We start by loading the credit card transaction dataset and performing exploratory data analysis (EDA). This includes examining the first few rows of the dataset, summary statistics, and data types to understand its structure. Visualizations such as correlation heatmaps and histograms are used to gain insights into the data distribution and relationships between features.

Data Augmentation with Geometric Distribution Mask

To enhance the robustness and generalization capabilities of our anomaly detection model, we employed a data augmentation technique using a geometric distribution mask. Specifically, we implemented a function that generates a mask based on a geometric distribution with a parameter (p) set to 0.5. This mask determines which elements of the original dataset are retained and which are set to zero, simulating the presence of missing or noisy data. By applying this mask to the training data, we create a new dataset where certain values are randomly zeroed out, encouraging the model to learn more robust representations that are less sensitive to individual data points. This approach helps the model to better capture the underlying patterns and dependencies within the multivariate time series data, thereby improving its ability to detect anomalies in diverse and noisy real-world scenarios. The masked data generated from this process was then used to train our transformer-based autoencoder within the anomaly detection framework.

Model Architecture:

We propose a deep learning-based architecture for anomaly detection comprising three main components:

1. **Transformer Autoencoder:** Utilizes a transformer-based encoder-decoder architecture to learn a low-dimensional representation of the input data and reconstruct it with minimal loss.
2. **Discriminator:** A binary classifier that distinguishes between normal and anomalous transactions.
3. **Contrastive Learning Model:** A self-supervised learning approach that maximizes similarity between similar instances and minimizes it between dissimilar ones, aiding in learning robust representations.

Model Training:

The combined model is trained in an end-to-end manner, optimizing the autoencoder's reconstruction loss, the discriminator's binary cross-entropy loss, and the contrastive learning model's mean squared error loss. We use the Adam optimizer and specify hyperparameters such as the number of epochs and batch size.

Evaluation Metrics:

To evaluate the performance of our anomaly detection methods, we consider several metrics including the F1 score and execution time. The F1 score provides a balance between precision and recall, while execution time reflects the computational efficiency of each method.

Comparison with Traditional Methods:

We compare the performance of our deep learning-based model with traditional anomaly detection methods, including PCA reconstruction error, OneClassSVM, Isolation Forest, Local Outlier Factor, DBSCAN, and Dynamic Time Warping. Each method is evaluated based on its F1 score and execution time, providing insights into their effectiveness and computational complexity.

Results and Discussion:

Our experimental results demonstrate that the proposed deep learning-based approach outperforms traditional methods in terms of anomaly detection accuracy. The Transformer Autoencoder achieves superior F1 scores compared to other techniques, indicating its effectiveness in capturing complex patterns in credit card transactions. However, it also requires more computational resources due to its deep architecture and parameter-intensive nature.

Conclusion:

In conclusion, we have presented a comprehensive approach for anomaly detection in credit card transactions using both traditional and deep learning-based methods. Our results highlight the efficacy of deep learning techniques in capturing intricate patterns and detecting fraudulent behavior. Future work may involve further optimization of the deep learning model and exploration of ensemble techniques for enhanced performance.