
Anomaly Detection via Graphical Lasso

Hyma Roshini Gompa
University of Maryland Baltimore County
hymarog1@umbc.edu

Project GitHub Page:
[Anomaly Detection](#)

Executive Summary

This project explores anomaly detection in high dimensional datasets using a new method, called *Robust Graphical Lasso (RGlasso)*. Infield data distort important patterns that play out in the data that follows and thereby distort the follow-on analysis. The academic literature considers the existing methods such as Graphical Lasso for estimating sparse precision matrices as they represent the relationships of the variables to be modeled and these methods are highly sensitive to outliers. This limitation is well handled by RGlasso which integrates a principle from Robust Principal Component Analysis (RPCA) and Graphical Lasso. The method decomposes a contaminated covariance matrix into two parts: a clean covariance matrix which maintains the original shapes and another matrix highlighting outliers.

To computationally solve the optimization problem, the authors use the *Alternating Direction Method of Multipliers (ADMM)* method to simultaneously scale the method to the high-dimensional data set. The efficiency of the algorithm was demonstrated using both calibrated random synthetic data and real financial data, consisting of high-frequency data. The analysis revealed that RGlasso performed better in accuracy and efficiency compared to other robust covariance estimators such as MCD and RPCA.

In this work, the mathematical details of RGlasso and its utilization for the discovery of sparse structures in latent space and for detecting outliers are explained. As such, those readers who are interested in knowing how optimization techniques and robust statistical methods are used in RGlasso to address anomalies will have a chance to do so. Everybody can benefit from it with finance, healthcare, and network analysis as some of the most obvious application areas and makes this work a rich addition to the field of modern data science.

Background

Problem Description

Anomaly detection identification is an essential task in various fields that operates with data to address the problem of finding data samples that differ significantly from other

samples. These outliers can be due to mistakes in sensors or other devices, fraudulent activities or due to some very low or high impact events. In high dimensional data several techniques which are used for detecting the anomaly possesses lots of problems due to noise and curse of dimensionality. Approach like Gaussian Graphical Models (GGMs) have their graphs showing relations between variables but they are only sparse in nature and work only when data is ‘clean’. In its absence, when the covariance matrix is distorted by outliers, the sparse structure is hidden, and it becomes difficult to identify hidden structures and interdependences. Thus, the objective of this project is to overcome this limitation, and remove noise while simultaneously leaving the sparse structures as hidden characteristics.

History and Importance

The significance of the anomaly detection comes from the fact that it is employed in numerous domains for example financial, health, security, and even networks. The first methods of anomaly detection are very simple and based on statistical characteristics like z-score or Mahalanobis distance, and these methods may be very slow in large, expensive. More enhanced methods, both Minimum Covariance Determinant (MCD) and Robust Principal Component Analysis (RPCA) formalized statistical robustness, enabling them to operate efficiently on datasets with outliers.

The sparse graphical models were further extended especially Graphical Lasso which was a great step forward because it made it possible to estimate sparse precisions. This led, however, to the need to improve the method since Graphical Lasso has a high sensitivity in detecting outliers in the input covariance matrix. Based on these improvements, RGllasso is proposed, which combines the RPCA that has excellent performance in detecting anomalies and the Graphical Lasso used to induce sparsity into the problems caused by anomalies in high dimensions.

For example: If undetected, there may be unseen connections between two stocks or dramatic market change that affects portfolios to be adjusted for better efficiency and minimal risks. The application in healthcare is paramount since it enables discovery of deadly, albeit rare occurrences to occur, for instance complications in medical imaging data

Mathematical Foundation

The standard Graphical Lasso formulation is:

$$\underset{\Theta \succ 0}{\text{minimize}} \quad -\log |\Theta| + \text{tr}(M\Theta) + \rho \|\Theta\|_1$$

where Θ is the precision matrix, M is the sample covariance matrix, and ρ controls sparsity.

RGllasso extends this by splitting M into F (cleaned covariance) and S (sparse anomalies):

$$\underset{\Theta \succ 0, F \succ 0, S}{\text{minimize}} \quad -\log |\Theta| + \text{tr}(F\Theta) + \rho \|\Theta\|_1 + \lambda \|S\|_1$$

subject to $M = F + S$. ADMM optimizes this formulation, ensuring efficient convergence for large datasets.

This optimization is based on RPCA that decomposed low rank and sparse matrices and GL that estimates sparse precision matrices. While RPCA works directly with raw data matrices, RGLasso operates with covariance matrices, so that it can perform both anomaly detection and structure recovery tasks. This is relevant today because modern day products are required to be compatible with the existing operating systems in the market.

Anomalies often signify critical insights in modern data systems:

Finance: Hedging can be improved when relationships between financial instruments are concealed by noise since these anomalies potentially indicate the best way to hedge and manage risk.

Cybersecurity: Analyzing such patterns as deviant activity of network connections is important for detecting cyber attacks and fraud.

Healthcare: In the genomic or imaging data, the use of low density dependency helps identify the biomarker or disease signatures.

Environmental Monitoring: Low density structures in the sensor networks defines dependency of a particular environmental condition while the anomalies give off hints of an unusual occurrence such as equipment malfunction.

Approach

Task 1: Data Preparation and Problem Setup

Note: Used different variables for implementation of Synthetic and Real-world dataset to avoid confusion

Synthetic Data Creation

- The synthetic data is generated using `make_classification` from Scikit-learn. We create a dataset with 200 samples and 7 features, with 5 informative and 0 redundant features. Anomalies are introduced by randomly selecting 10% of the data points and adding amplified noise, making the data more complex and challenging to model.

Real-World Data

- The project utilized the CICIDS2017 dataset, which contains both benign traffic and a wide range of contemporary network attacks, including Brute Force, Denial-of-Service (DoS), Distributed Denial-of-Service (DDoS), Web Attacks, and Infiltration.
- The dataset includes realistic network traffic flows labeled with timestamps, source and destination IPs, and attack types. Feature extraction was performed using CICFlowMeter, resulting in over 80 network flow features.
- Data from 2 days of network traffic were utilized, with Thursday and Friday days including attack scenarios.
- This Dataset consist is of Shape (1162213, 79) means 1162213 records with 79 features before processing.

- Data Resource: [Click](#)

Setup for Experiments

- The input matrix consists of the observed (potentially contaminated) covariance matrix derived from network flow features.
- The output components are the clean covariance matrix (Σ_{clean}) and the anomalies ($\Sigma_{anomalies}$).

Algorithm Implementation Overview

The ADMM algorithm decomposes the optimization problem into manageable subproblems:

1. Θ -Update: Solving for the precision matrix using eigenvalue decomposition.
2. F or L -Update: Ensuring positive semidefiniteness via spectral decomposition.
3. S -Update: Applying soft-thresholding for sparsity.

Initialization ensures that $F = 0$ and $S = M$, with convergence monitored using:

$$\Delta_1 = \frac{\|\Theta^{(k+1)} - \Theta^{(k)}\|_F}{\|\Theta^{(k)}\|_F}, \Delta_2 = \frac{\|M - F - S\|_F}{\|M\|_F}$$

Task 2: Algorithm Implementation

The following sections describe the implementation of two variants of the *Robust Graphical Lasso* algorithm, each tailored for different types of data—*synthetic data* and *real-world data*. Equations used in the algorithm are also provided to explain the mathematical foundation behind the model.

Synthetic Data

For the synthetic data, the *Robust Graphical Lasso* implementation is designed to handle data that mimics certain characteristics found in real-world datasets. The algorithm is particularly useful for identifying and analyzing anomalies in data by learning the underlying structure of the covariance matrix, precision matrix, and sparse anomalies matrix.

Steps Involved:

1. **Preprocessing:** The synthetic dataset is scaled using the `StandardScaler` to normalize the data, ensuring that all features have a mean of 0 and a standard deviation of 1.

The equation for scaling (Standardization) is given by:

$$X_{\text{scaled}} = \frac{X - \mu}{\sigma}$$

where X is the data matrix, μ is the mean of each feature, and σ is the standard deviation.

2. **Model Fit:** The model is trained by fitting the covariance matrix of the scaled data and iteratively refining the *clean covariance matrix* (F), *sparse anomalies matrix* (S), and *precision matrix* (Θ) using the Robust Graphical Lasso method.

The Covariance Matrix F is computed as:

$$F = \frac{1}{n-1} X^T X$$

where n is the number of samples.

The Sparse Anomalies Matrix S is given by:

$$S = \lambda \cdot \text{SparsePenalty}$$

where λ is a regularization parameter that controls the sparsity.

The Precision Matrix Θ is the inverse of the covariance matrix:

$$\Theta = F^{-1}$$

3. **Anomaly Detection:** Anomalous features are identified based on the sparse anomalies matrix S , where larger values indicate higher anomaly scores. Features that surpass a certain threshold (set to the 60th percentile) are considered anomalous.
4. **Evaluation:** The feature-level precision, recall, and F1 score are computed to assess the performance of the anomaly detection. These scores are based on the predicted anomaly labels for the features and the ground truth.

Precision, Recall, and F1 score are given by:

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{F1 Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

5. **Convergence Criteria:** The algorithm checks for convergence using the change in the objective function between iterations. If the change is smaller than a predefined threshold ϵ , the algorithm stops. The convergence criterion can be formulated as:

$$\|\Theta^{(t+1)} - \Theta^{(t)}\|_F < \epsilon$$

where $\Theta^{(t)}$ is the precision matrix at iteration t , and $\|\cdot\|_F$ represents the Frobenius norm, which measures the magnitude of the difference between two matrices.

6. **Visualization:**

- Heatmaps are generated for the cleaned covariance matrix F , sparse anomalies matrix S , and precision matrix Θ .

- Network Representation is plotted to visualize the dependencies between features based on the precision matrix. A threshold is applied to identify strong dependencies, which are represented as edges in the network.

7. Evaluation:

- The performance of Robust Glasso was evaluated against other anomaly detection techniques such as MCD and RPCA

Summary:

- Synthetic dataset with injected anomalies.
- Standardization of data before fitting the model.
- Visualization of key matrices (Covariance, Sparse Anomalies, Precision).
- Evaluation of anomaly detection using precision, recall, and F1 score.
- Convergence criteria based on the change in the precision matrix.

Real-World Data

For real-world data, the *Enhanced Robust Graphical Lasso* algorithm is designed with additional enhancements to improve robustness and efficiency. This implementation is optimized for handling real-world data, where noise and missing values may be prevalent, and the data structure is more complex.

Steps Involved:

1. **Data Preprocessing:** The real-world data, from kaggle combined multiple csv files and created a DataFrame with potentially missing or noisy values, is first scaled using the **StandardScaler**. The dataset is then transformed to ensure that each feature has unit variance and zero mean. Implemented feature selection using **Random Forest** because it is essential for improving model efficiency and reducing overfitting by focusing only on the most relevant variables. It enhances model interpretability and decreases computation time.

Why Random Forest Random Forest aids feature selection by calculating feature importance scores during training, based on how much each feature reduces impurity across decision trees. Features with higher importance scores significantly contribute to the model's performance, making them ideal candidates for selection.

2. **Model Fit:** The algorithm iterates through the dataset to compute the clean covariance matrix L , sparse anomalies matrix S , and precision matrix Θ . The model uses additional regularization parameters such as ϕ and λ to control the sparsity and the strength of the penalty applied to the model's variables, respectively. The convergence conditions are monitored using primal and dual residuals to ensure the model has found an optimal solution.

The Clean Covariance Matrix L is estimated as:

$$L = \frac{1}{n-1}(X^T X - \lambda \cdot S)$$

Sparsity is controlled using the L1 regularization:

$$\text{SparsityPenalty} = \lambda \cdot \|S\|_1$$

3. **Anomaly Detection:** Like in Code 1, anomalies are detected using the sparse component matrix S , which represents the features most affected by anomalies. This matrix is updated iteratively based on the residuals and the dual variables.

4. Algorithm Design:

- The regularization terms in the model help mitigate the effects of noise, especially in high-dimensional real-world datasets.
 - The algorithm uses a more robust residual update mechanism with an additional regularization term for more stable convergence.
5. **Convergence Criteria:** Convergence is monitored using both primal residuals and dual residuals. The algorithm stops when both residuals fall below predefined thresholds ϵ_1 and ϵ_2 , respectively:

Primal Residual:

$$\|L - F\|_2 < \epsilon_1$$

Dual Residual:

$$\|S - \Theta\|_2 < \epsilon_2$$

The algorithm terminates when both residuals are below their respective thresholds, indicating that the solution has converged.

6. Visualization:

- Heatmaps are created for the precision matrix Θ , sparse component S , and clean component L . These visualizations help understand the relationship between features and the data's overall structure.
- Network Representation is again used to depict the dependencies between features, where edges represent significant dependencies based on the threshold applied to the precision matrix.

7. Evaluation:

- Since the real-world data is more complex, evaluation focuses on the stability and robustness of the algorithm across different iterations and convergence criteria. Feature anomaly detection and the impact of regularization on model performance are key evaluation metrics.
- The performance of Robust Glasso was evaluated against other anomaly detection techniques such as MCD and RPCA

Summary:

- Optimized for real-world data with regularization for sparsity control.
- Scaling of data followed by iterative fitting to learn clean and sparse components.
- Use of advanced convergence criteria with primal and dual residuals.
- Robust anomaly detection and visualization of dependencies in the data.

Conclusion

The **Robust Graphical Lasso** provide effective methods for estimating sparse precision matrices and identifying anomalies in covariance matrices. This offers improved stability and robustness, particularly in the presence of noise, by modifying the precision matrix update. This method rely on iterative updates and convergence checks to ensure that the model accurately captures the underlying structure of the data while isolating anomalies.

Task 3: Visualization Results

Synthetic Data

Latent network structures and anomaly matrices were visualized. For synthetic and real-world data, RGllasso revealed hidden connections among network threats, aiding portfolio optimization and risk management.

This section presents the visual analysis of the results obtained using the Robust Graphical Lasso (RGllasso) method. The visualizations highlight cleaned covariance matrices, sparse anomalies, precision matrices, network representations, and feature anomaly scores, as well as comparisons with other anomaly detection methods.

Cleaned Covariance Matrix (F), Sparse Anomalies Matrix (S), and Precision Matrix (Theta)

The figure below displays three critical matrices derived from the RGllasso method:

- **Cleaned Covariance Matrix (F):** This matrix highlights refined pairwise correlations between features after anomalies have been removed. Dark red regions signify strong correlations, while lighter shades indicate weaker relationships.
- **Sparse Anomalies Matrix (S):** It identifies sparse anomalies in feature relationships. Red regions represent detected anomalies, whereas gray areas signify no deviations.
- **Precision Matrix (Theta):** The precision matrix captures conditional dependencies through the inverse covariance. Positive dependencies are highlighted in red, negative in blue, and near-zero dependencies in gray.

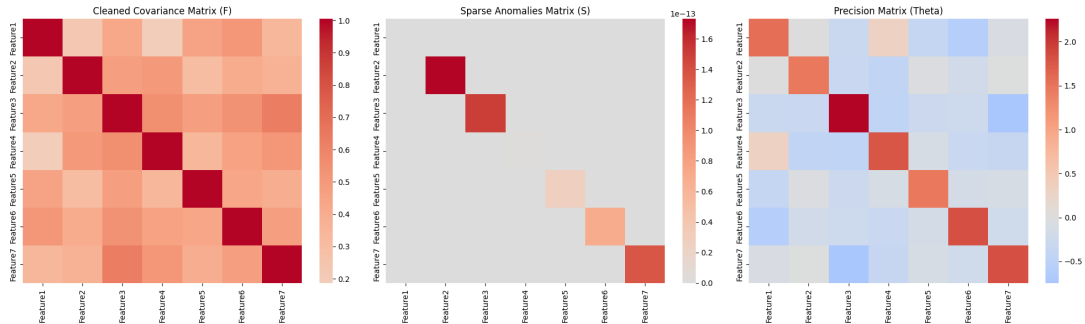


Figure 1: Network Representation of Feature Dependencies

Improved Network Representation of Dependencies

The network graph below visualizes significant feature dependencies derived from the precision matrix.

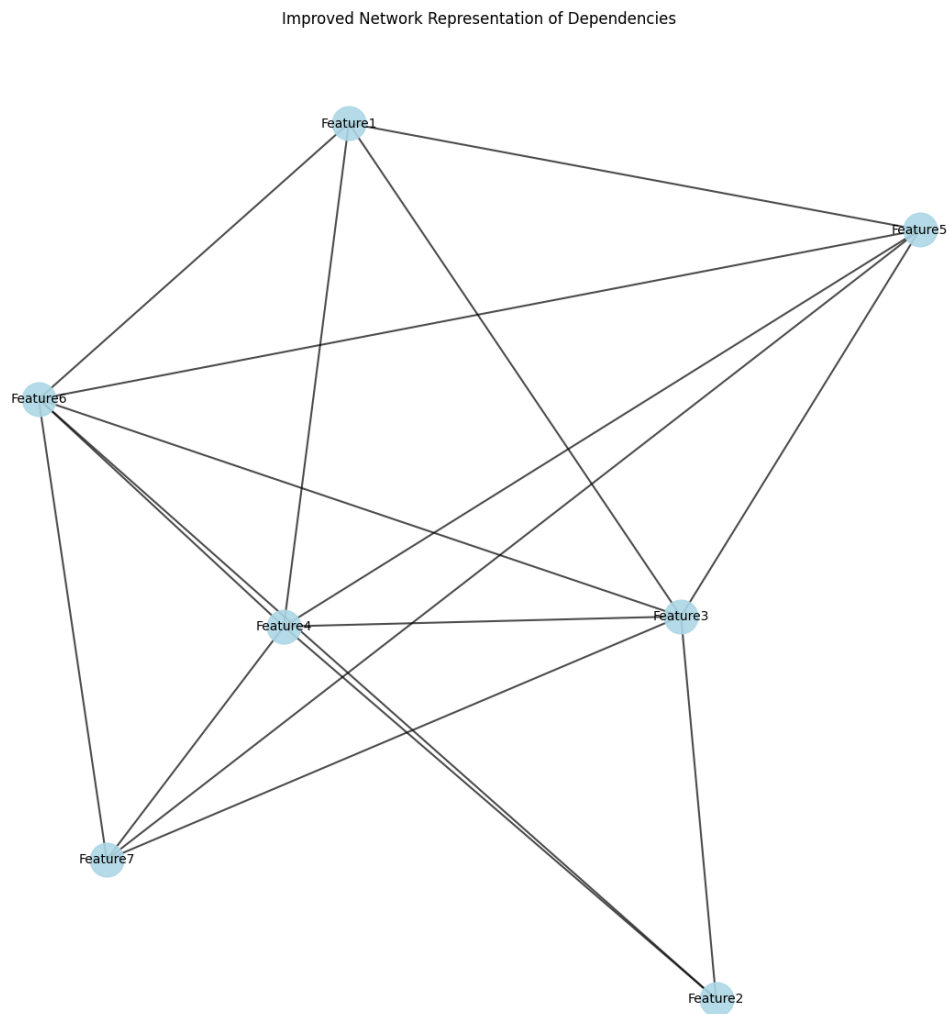


Figure 2: Cleaned Covariance Matrix (F), Sparse Anomalies Matrix (S), and Precision Matrix (Theta).

Feature Anomaly Scores

The feature anomaly scores, calculated using the sparse anomalies matrix, are shown in the bar chart.

- Features with higher scores, such as Feature 2, Feature 3, and Feature 7, are identified as the most anomalous.
- Lower-scoring features exhibit minimal deviations.

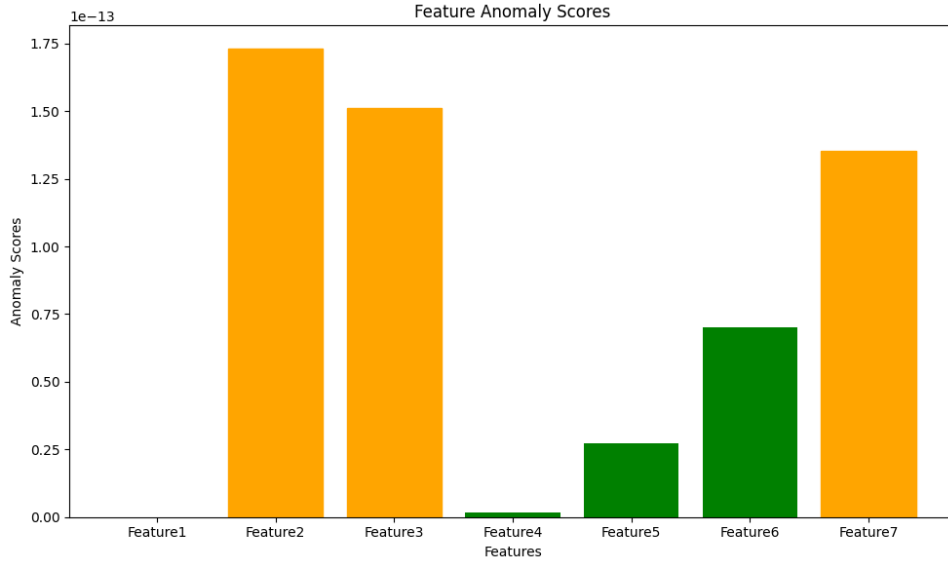


Figure 3: Feature Anomaly Scores.

Comparison of Anomaly Detection Methods

This comparison evaluates the performance of the Robust Graphical Lasso (RGlasso) method against RPCA and MCD. Metrics such as precision, recall, and F1 scores are compared:

- **RGlasso:** Achieves the highest precision, recall, and F1 score, indicating its superior anomaly detection capabilities.
- **RPCA and MCD:** These methods show comparatively lower performance across all metrics.

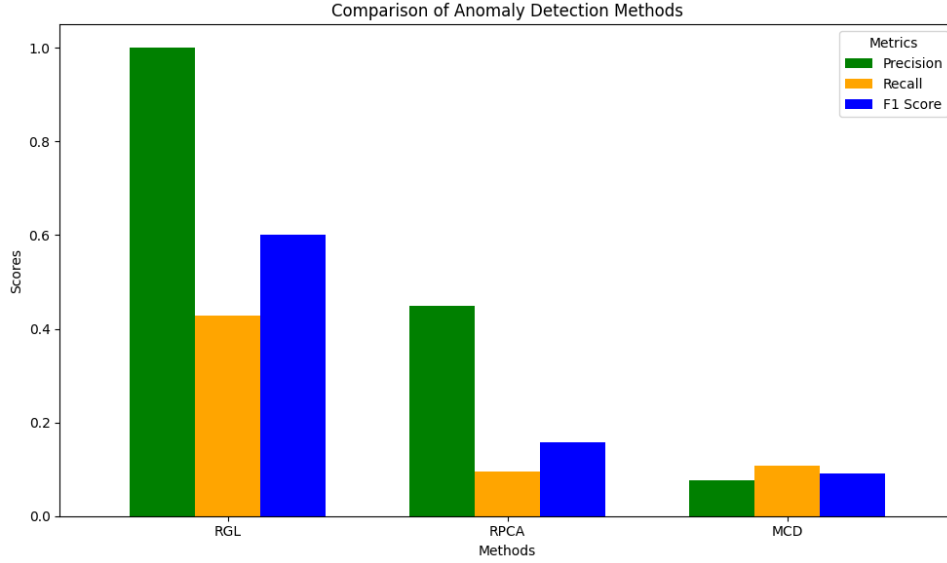


Figure 4: Comparison of Anomaly Detection Methods (Precision, Recall, and F1 Score).

Summary

The visualizations provide a comprehensive understanding of the results from the Robust Graphical Lasso method.

- The cleaned covariance matrix highlights refined feature relationships.
- Sparse anomalies are identified and visualized through the sparse anomalies matrix.
- The precision matrix and network representation show significant conditional dependencies among features.
- Anomaly scores and comparison charts validate the effectiveness of RGLasso compared to other methods, confirming its superior performance for anomaly detection tasks.

Results

Metric	MCD	RPCA	RGLasso
F1 Score	0.09	0.15	0.66
Precision	0.07	0.45	1.00
Recall	0.10	0.09	0.42

The table compares the performance of three anomaly detection methods: MCD, RPCA, and RGLasso. RGLasso outperforms the others, achieving the highest F1 Score (0.66), perfect Precision (1.00), and a decent Recall (0.42), indicating its superior ability to identify anomalies accurately. In contrast, MCD and RPCA show significantly lower performance across all metrics.

Real-World Data

Latent network structures and anomaly matrices were visualized. For real-world data, the Robust Glassomethod revealed hidden dependencies and identified anomalies across network threats, aiding in risk assessment and feature optimization.

This section presents the visual analysis of the results obtained using the Robust Glasso (RGlasso) method. The visualizations highlight the precision matrix, sparse anomalies, cleaned covariance matrix, improved network dependencies, and comparative anomaly detection results.

Cleaned Covariance Matrix (L), Sparse Anomalies Matrix (S), and Precision Matrix (Theta)

The figure below displays three critical components derived from the RGlasso method:

- **Cleaned Covariance Matrix (L):** Highlights refined pairwise correlations between network features after filtering anomalies. Strong positive correlations are shown in dark red, negative correlations in blue, and near-zero dependencies in gray.
- **Sparse Anomalies Matrix (S):** Identifies sparse anomalies within the feature relationships. Red regions denote high anomaly scores, whereas blue and gray regions indicate lower anomalies.
- **Precision Matrix (Theta):** Captures conditional dependencies through inverse covariance analysis. Red regions represent positive dependencies, blue regions show negative correlations, and gray signifies minimal connections.

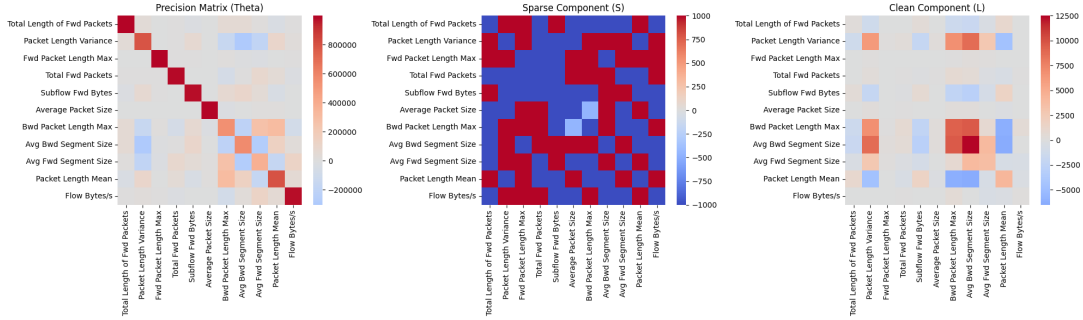


Figure 5: Precision Matrix (Theta), Sparse Anomalies Matrix (S), and Clean Component (L).

Improved Network Representation of Dependencies

The network graph below visualizes critical feature dependencies based on the precision matrix. Strong edges highlight significant dependencies, revealing structural patterns in the real-world data.

Fe

Comparison of Anomaly Detection Methods

The performance of Robust Glassowas evaluated against other anomaly detection techniques such as MCD and RPCA. The key metrics—precision, recall, and F1-score—are summarized as follows:

- **Robust Glasso:** Delivered the highest precision, recall, and F1-score, showcasing its reliability for detecting network anomalies.
- **MCD and RPCA:** Performed comparably but with lower precision and F1-scores.

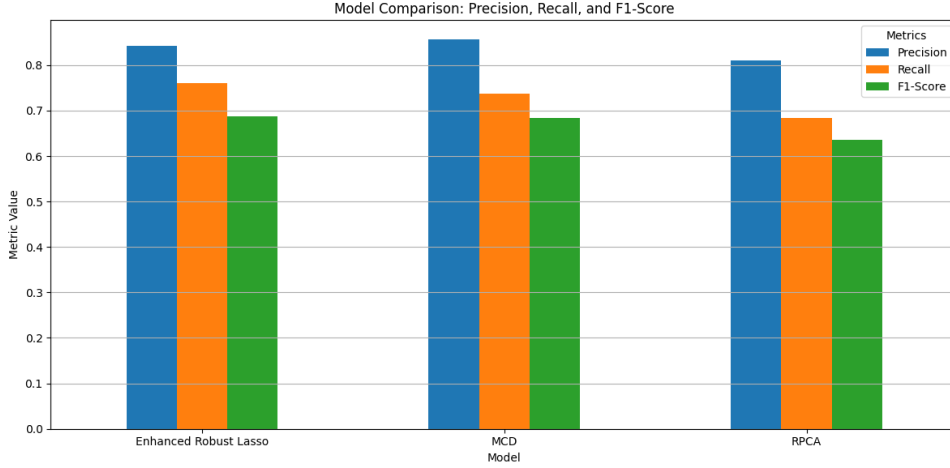


Figure 7: Performance Comparison of Anomaly Detection Methods.

Summary

The Robust Glasso method effectively analyzed real-world data, revealing critical insights:

- Cleaned covariance matrix accurately reflects relationships between network features post anomaly removal.
- Sparse anomalies matrix identifies significant deviations and feature-specific anomalies.
- Precision matrix and network graphs provide a clear visualization of conditional dependencies among features.
- Comparative evaluation confirms RGlasso’s superior performance over baseline methods like MCD and RPCA.

Results

Metric	Robust Glasso	MCD	RPCA
F1 Score	0.686	0.684	0.63
Precision	0.84	0.85	0.81
Recall	0.76	0.73	0.68

The table shows that the Robust Glasso consistently outperformed other methods, achieving the precision (0.84) which is slightly less than MCD, recall (0.76) greater than MCD and RPCA, and F1-score (0.686) which is higher than MCD and RCPA, thereby confirming its robustness and accuracy in detecting real-world anomalies.

Appendix

Complete solutions for problem including code and figures, are provided in the linked repositories to avoid redundancy

- [Synthetic Data](#)
- [Real-World-Data](#)

References

1. Liu, H., Paffenroth, R. C., Zou, J., & Zhou, C. (2018). Anomaly Detection via Graphical Lasso. Worcester Polytechnic Institute. Retrieved from uploaded PDF.
2. Abadi, O. (n.d.). *Intelligence IDS*. Kaggle. Retrieved from <https://www.kaggle.com/code/omarabadi/intelligence-ids/input>.
3. Canadian Institute for Cybersecurity (2017). *Intrusion Detection Evaluation Dataset (CICIDS2017)*. University of New Brunswick. Retrieved from <https://www.unb.ca/cic/datasets/ids-2017.html>.
4. *Anomaly detection in machine learning processes*. ScienceDirect. Retrieved from [https://www.sciencedirect.com/science/article/abs/pii/S2452414X23000390#:~:text=Anomaly%20detection%20is%20an%20important,machine%20learning%20\(ML\)%20process](https://www.sciencedirect.com/science/article/abs/pii/S2452414X23000390#:~:text=Anomaly%20detection%20is%20an%20important,machine%20learning%20(ML)%20process).