# Unsupervised Machine Learning

# Anomaly Detection in Network Traffic

**Presented by:**

**Sancia Fernandes (A012)**

**Sherin Ouseph (A017)**

**PROBLEM STATEMENT:**

*Develop an anomaly detection system for network traffic using unsupervised machine learning techniques. The system should be capable of identifying unusual patterns and potential threats within the network traffic data without relying on labeled examples. The primary objective is to design a robust solution that can effectively detect anomalies such as intrusions or suspicious activities within the network.*

**OBJECTIVE:**

*The objective of this project is to build an unsupervised anomaly detection system tailored for network traffic analysis, with the primary aim of pinpointing irregularities indicative of potential security breaches or anomalous behavior within the network. The target variable encompasses instances classified as anomalies or normal network traffic patterns. The objective is to develop a solution capable of autonomously identifying and flagging suspicious activities without the need for prior labeling, thereby bolstering network security measures and safeguarding against potential threats.*

**RESEARCH METHODOLOGY:**

**Data Collection:** Acquire network traffic data from relevant sources, ensuring it covers a diverse range of network activities and anomalies.

**Data Preprocessing:**

Handle missing values: Impute missing values or remove instances with missing data.
Encode categorical variables: Convert categorical features into numerical representations using techniques like one-hot encoding.
Scale numerical features: Normalize or standardize numerical features to ensure uniformity in scale.

**Exploratory Data Analysis (EDA):**

Analyze correlations: Examine pairwise correlations between features to identify potential relationships.
Identify outliers: Use visualization techniques such as scatter plots or box plots to detect outliers in the data.

**Unsupervised Learning Algorithms:**

Implement KMeans: Partition the data into clusters based on similarity, experimenting with different numbers of clusters.
Apply DBSCAN: Identify dense regions of data points as clusters while labeling outliers as anomalies.
Utilize Isolation Forest: Train an ensemble of decision trees to isolate anomalies based on their low prevalence in the data.
Employ One-Class SVM: Train a support vector machine model to identify outliers as instances lying far from the majority of data points.

**Dimensionality Reduction:**

Use Principal Component Analysis (PCA): Reduce the dimensionality of the data while preserving as much variance as possible. Visualize the reduced-dimensional data to aid in clustering and anomaly detection.

**Evaluation and Validation:**

Assess clustering quality: Calculate silhouette score, adjusted Rand index, and other clustering metrics to evaluate the quality of clusters formed by different algorithms.
Evaluate anomaly detection models: Measure the accuracy, precision, recall, and F1-score of anomaly detection models using appropriate evaluation metrics.
Validate model performance: Utilize cross-validation techniques to validate the robustness and generalization ability of the anomaly detection system.
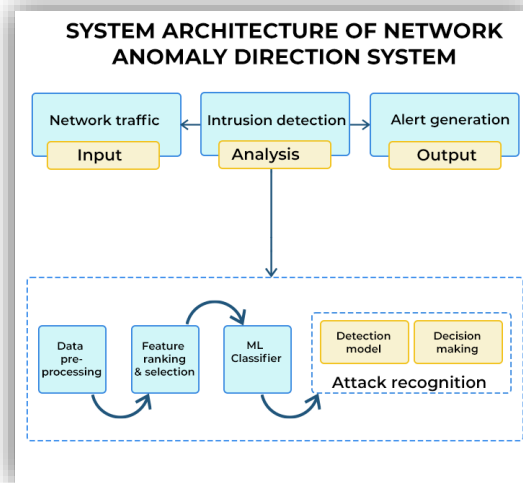Documentation and Reporting:
Document all preprocessing steps, algorithm implementations, and evaluation results.
Prepare a comprehensive report detailing the methodology applied, findings, and recommendations for further improvements.
Present insights gained from the analysis and discuss the implications for network security and anomaly detection.

## INTRODUCTION:



In today's interconnected world, network security is of paramount importance to safeguard sensitive information and ensure the uninterrupted operation of critical systems. With the proliferation of cyber threats and sophisticated attack techniques, organizations face the constant challenge of detecting and mitigating potential security breaches within their network infrastructure. Traditional security measures, such as firewalls and intrusion detection systems, are often insufficient in detecting novel and stealthy threats.

To address this challenge, the field of anomaly detection has emerged as a promising approach to identifying suspicious activities and anomalies within network traffic data. Anomaly detection systems aim to differentiate between normal network behavior and abnormal patterns that may indicate malicious activity or system failures. By leveraging unsupervised machine learning techniques, these systems can autonomously identify anomalies without the need for labeled examples, thereby providing a proactive means of enhancing network security.

In this project, we propose to develop a robust anomaly detection system tailored for network traffic analysis. The objective is to design a solution capable of accurately identifying and flagging anomalous patterns within network data, including potential intrusions, suspicious activities, or abnormal behaviors. By leveraging a combination of data preprocessing, exploratory analysis, unsupervised learning algorithms, and dimensionality reduction techniques, we aim to build a comprehensive anomaly detection pipeline that can effectively detect and respond to security threats in real-time.

The project encompasses a systematic methodology, starting from data collection and preprocessing to model implementation, evaluation, and documentation. Through thorough analysis and experimentation, we seek to develop an anomaly detection system that not only achieves high detection accuracy but also offers scalability and robustness in real-world network environments. By providing timely detection and response to anomalous behavior, the proposed system aims to bolster network security measures and safeguard against potential cyber threats.

## DATA COLLECTION:

For this project, the network traffic data was obtained from Kaggle, a popular platform for sharing datasets and machine learning resources. The dataset contains a comprehensive collection of network traffic records captured from various sources, providing insights into different types of network activities and potential anomalies.

The data collection process involved accessing the dataset from its source on Kaggle and downloading it for further analysis. The dataset comprises structured data in a CSV (Comma Separated Values) format, making it suitable for manipulation and analysis using tools like Python and pandas.

Upon downloading the dataset, an initial inspection was conducted to understand its structure, features, and size. The dataset includes information such as protocol types, service types, flags, and attack types, among others. Each record represents a network traffic instance captured at a specific point in time, with attributes describing various aspects of the communication.

It is important to note that the dataset may contain both normal and anomalous network traffic instances, providing a diverse and realistic representation of network behavior. This diversity is crucial for training and evaluating anomaly detection models effectively.

Furthermore, the dataset require preprocessing steps such as encoding categorical variables, and scaling numerical features before it can be used for analysis and model training. These preprocessing steps ensure the data's consistency and prepare it for subsequent stages of the project, such as exploratory data analysis and model implementation.

## DIMENSIONALITY REDUCTION :

In this unsupervised machine learning (UML) project, Principal Component Analysis (PCA) is used for dimensionality reduction. PCA is applied to reduce the number of features (columns) in the dataset while retaining the most important information. Here's a breakdown of PCA and its application in this project:

Principal Component Analysis (PCA):

PCA is a technique used to transform high-dimensional data into a lower-dimensional space while preserving the maximum variance in the data. It achieves this by identifying the principal components, which are new variables that are linear combinations of the original variables. These components are orthogonal (uncorrelated) to each other.
The first principal component (PC1) explains the maximum variance in the data, followed by PC2, PC3, and so on. PCA orders the principal components by the amount of variance they explain, allowing us to select a subset of components that capture most of the variability in the data.
It is commonly used for dimensionality reduction, visualization, noise reduction, and feature extraction.

Why PCA is Applied:

PCA is applied to address the issues related to high-dimensional data, such as the curse of dimensionality, computational complexity, and visualization challenges.

It helps in capturing the most important information in the data while discarding less relevant information, thus simplifying the subsequent analysis.

PCA aids in understanding the underlying structure of the data by visualizing it in a lower-dimensional space, facilitating better interpretation and decision-making.

By reducing the number of features, PCA can improve the efficiency and effectiveness of clustering algorithms, making them more scalable and accurate.

Application:

PCA is applied after encoding categorical variables using one-hot encoding and before performing clustering algorithms like KMeans and DBSCAN. The original dataset contain a large number of features after one-hot encoding categorical variables, which can lead to the curse of dimensionality and computational inefficiency, especially for clustering algorithms. PCA is used to reduce the dimensionality of the dataset from a high-dimensional space to a lower-dimensional space (in this case, from many one-hot encoded features to just two principal components). By reducing the dimensionality, PCA helps in visualizing the data in a two-dimensional space, making it easier to observe clusters and patterns. Additionally, reducing the dimensionality can improve the performance of clustering algorithms by reducing noise and computational complexity, as clustering tends to work better in lower-dimensional spaces.

## UNSUPERVISED MACHINE LEARNING ALGORITHMS:

### 1. KMEANS:

*KMeans is a popular clustering algorithm used in unsupervised machine learning for partitioning a dataset into distinct clusters based on similarity. The algorithm aims to minimize the within-cluster variance, where each data point is assigned to the cluster with the nearest centroid. KMeans works iteratively to optimize cluster centroids until convergence, typically defined by a predefined number of iterations or when centroids no longer change significantly.*

In this project, KMeans is utilized as one of the clustering algorithms for anomaly detection in network traffic data. The algorithm partitions the dataset into clusters based on the similarity of network traffic instances. By analyzing the characteristics of each cluster, anomalies or outliers can be identified as data points lying far from the centroids of normal clusters.

KMeans does not require labeled data for training, making it suitable for detecting anomalies in network traffic without prior knowledge of specific attack types or anomalies. It autonomously identifies patterns and clusters within the data based on similarities in network behavior.

Scalability: KMeans is computationally efficient and scalable to large datasets, making it suitable for analyzing extensive network traffic records captured over extended periods. It can handle high-dimensional data efficiently, making it suitable for analyzing network traffic with multiple features.

Interpretability: The clusters formed by KMeans provide interpretable insights into the underlying structure of the network traffic data. By examining the characteristics of each cluster, anomalies or outliers can be identified based on their deviation from normal cluster patterns.
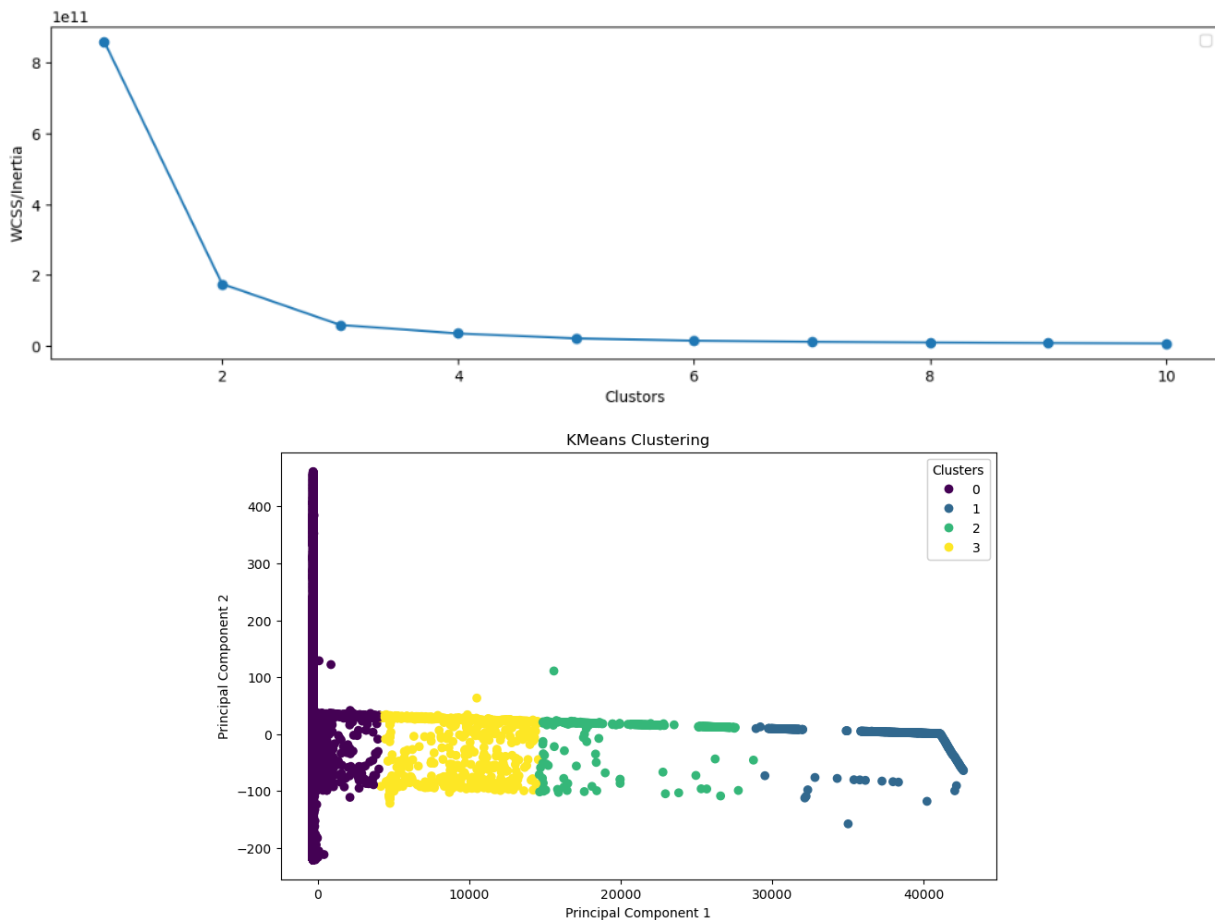
Flexibility: KMeans allows for the specification of the number of clusters (K), providing flexibility in adjusting the granularity of anomaly detection. By experimenting with different values of K, the algorithm can adapt to varying levels of complexity and heterogeneity in network traffic data.

The elbow method is a technique used to determine the optimal number of clusters (K) for K-means clustering. The idea is to plot the within-cluster sum of squares (WCSS) as a function of K and look for a "knee" or an "elbow" in the plot, which indicates the point where increasing the number of clusters does not significantly improve the clustering performance.

In the context of anomaly detection in network traffic, K-means clustering can be used to identify clusters of "normal" network traffic and detect anomalies as data points that do not belong to any of the normal clusters. The elbow method can be used to determine the optimal number of clusters that best represent the normal traffic patterns.

Overall, K-means clustering and the elbow method can be effective techniques for anomaly detection in network traffic. However, it is important to note that these techniques may not always detect all types of anomalies, especially if they are rare or sophisticated. Therefore, it is recommended to use a combination of techniques and approaches for comprehensive anomaly detection in network traffic.

Overall, KMeans serves as a valuable tool for clustering network traffic data and identifying anomalies through unsupervised learning. Its simplicity, efficiency, and interpretability make it well-suited for anomaly detection tasks in network security applications.





## 2. DBSCAN (Density-Based Spatial Clustering of Applications with Noise):

*DBSCAN is a density-based clustering algorithm used in unsupervised machine learning for identifying clusters in a dataset with varying densities. Unlike KMeans, which requires the number of clusters to be specified in advance, DBSCAN can automatically discover clusters of arbitrary shapes and sizes.*

DBSCAN works by partitioning the dataset into three types of points:

Core Points: These are data points that have a specified minimum number of neighbors (defined by the parameters epsilon and min_samples) within a defined radius. Core points lie within dense regions of the dataset.

Border Points: These points are within the epsilon radius of a core point but do not have enough neighbors to be considered core points themselves. Border points are on the fringes of clusters.

Noise Points: Data points that do not belong to any cluster and do not have enough neighbors to be considered core points or border points. These points are considered outliers or noise.

DBSCAN is used in this project for anomaly detection in network traffic data for several reasons:
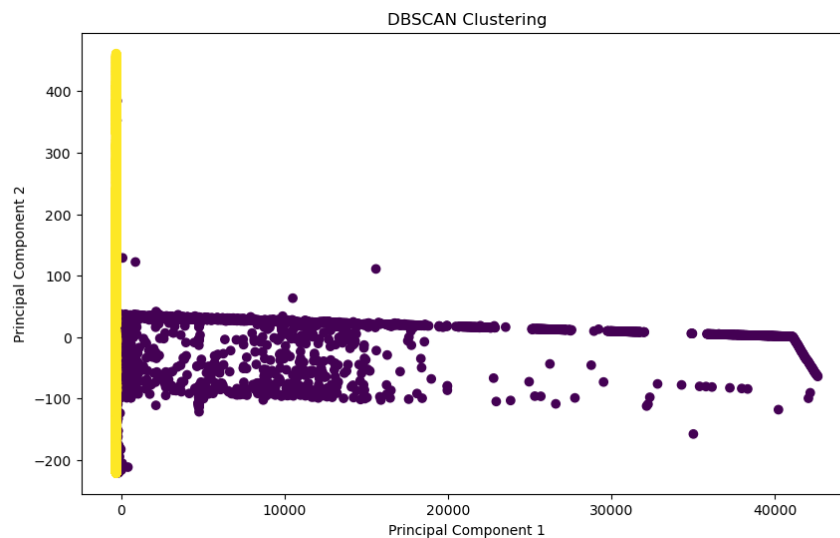
Flexibility in Cluster Shapes: DBSCAN is capable of identifying clusters of arbitrary shapes and sizes, making it suitable for detecting anomalies that may not conform to traditional cluster structures. This flexibility allows DBSCAN to capture complex patterns in network traffic data, including both global and local anomalies.

Robustness to Noise: DBSCAN is robust to noise and outliers in the dataset. It effectively distinguishes between true clusters and noise points, enabling the detection of anomalous network traffic instances that deviate significantly from the underlying cluster patterns.

Automatic Determination of Cluster Density: DBSCAN automatically determines cluster density based on the epsilon parameter, which defines the maximum distance between two points to be considered neighbors. This adaptive approach allows DBSCAN to adapt to variations in the density of network traffic data and identify clusters accordingly.

No Assumption of Cluster Shape: Unlike parametric clustering algorithms like KMeans, DBSCAN does not make assumptions about the shape or size of clusters. It can detect clusters of varying shapes, densities, and sizes, making it well-suited for analyzing heterogeneous network traffic data.

Overall, DBSCAN offers a powerful and versatile approach to anomaly detection in network traffic data. Its ability to identify clusters of arbitrary shapes, robustness to noise, and automatic determination of cluster density make it a valuable tool for detecting anomalies and identifying potential security threats within network traffic.



### 3. ONE-CLASS SVM:

*One-Class SVM is a variant of the traditional SVM algorithm used for anomaly detection in unsupervised learning scenarios where only one class of data is available for training. It learns a*

*decision boundary around the majority of the data points, considering them as normal instances, while treating instances lying outside this boundary as anomalies.*

*One-Class SVM works by finding the hyperplane that best separates the normal instances from the origin in feature space. The objective is to maximize the margin around the normal data points while minimizing the number of data points lying outside the margin. This approach effectively identifies anomalies as data points lying outside the learned boundary.*

One-Class SVM is used in this project for anomaly detection in network traffic data for several reasons:
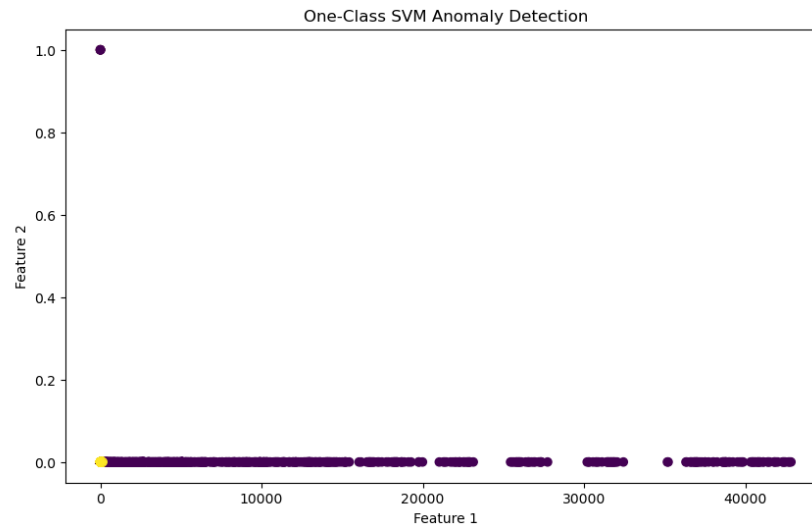
Unsupervised Anomaly Detection: One-Class SVM does not require labeled examples of anomalies for training. Instead, it learns the characteristics of normal network traffic instances and identifies anomalies as deviations from this normal behavior. This makes it suitable for detecting novel and previously unseen anomalies in network traffic data.

Robust to Outliers: One-Class SVM is robust to outliers and noise in the dataset. It learns a decision boundary around the majority of the data points, effectively ignoring outliers and focusing on capturing the underlying structure of the normal data distribution. This robustness is crucial for detecting anomalies in network traffic data, which may contain noisy or irregular instances.

Flexibility in Feature Space: One-Class SVM can operate in high-dimensional feature spaces, making it suitable for analyzing complex network traffic data with multiple features. It can capture nonlinear relationships and interactions between features, allowing for effective anomaly detection in high-dimensional data.

Scalability: One-Class SVM is computationally efficient and scalable to large datasets, making it suitable for analyzing extensive network traffic records captured over extended periods. It can handle high-dimensional data efficiently, making it well-suited for detecting anomalies in network traffic data with multiple features.

Overall, One-Class SVM offers a powerful and versatile approach to anomaly detection in network traffic data. Its ability to detect anomalies without the need for labeled examples, robustness to outliers, flexibility in feature space, and scalability make it a valuable tool for identifying potential security threats and anomalies within network traffic.

## 4. ISOLATION FOREST

*Isolation Forest is an unsupervised machine learning algorithm used for anomaly detection, particularly in scenarios where anomalies are expected to be rare and distinct from the majority of normal instances. It operates by isolating anomalies in the dataset by recursively partitioning the feature space into smaller subsets.*

*The key idea behind Isolation Forest is that anomalies are likely to be isolated in the feature space, meaning they require fewer splits to be separated from the majority of normal instances. In contrast, normal instances are expected to require more splits to isolate them, as they are densely packed together. By measuring the number of splits required to isolate each instance, anomalies can be identified as instances that are isolated with fewer splits.*

Isolation Forest is used in this project for anomaly detection in network traffic data for several reasons:

Efficient for High-Dimensional Data: Isolation Forest is efficient in handling high-dimensional data, making it suitable for analyzing network traffic data with multiple features. It can effectively identify anomalies in complex feature spaces with varying dimensions.
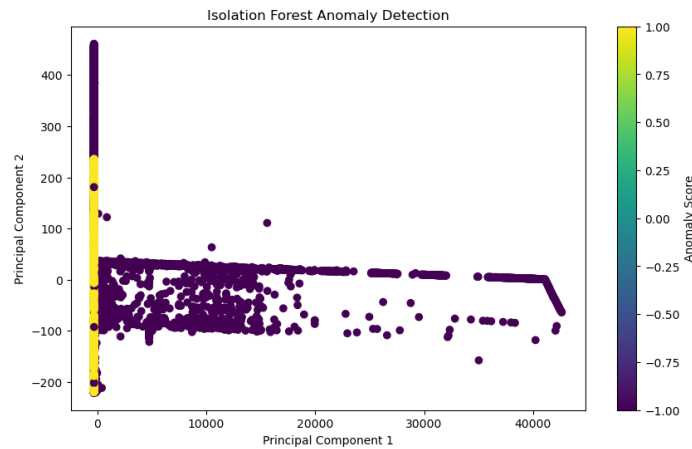
Scalability: Isolation Forest is computationally efficient and scalable to large datasets, making it suitable for analyzing extensive network traffic records captured over extended periods. It can handle large volumes of data efficiently, making it well-suited for real-time anomaly detection in network traffic.

Robust to Noise: Isolation Forest is robust to outliers and noise in the dataset. It isolates anomalies by focusing on the intrinsic structure of the data rather than the specific characteristics of outliers. This robustness allows it to effectively detect anomalies in noisy or irregular network traffic data.

No Assumptions about Data Distribution: Isolation Forest does not assume any specific data distribution or cluster shape. It can detect anomalies in datasets with arbitrary distributions and shapes, making it suitable for analyzing heterogeneous network traffic data with diverse patterns and behaviors.

Fast Training and Inference: Isolation Forest has a simple and fast training process, making it suitable for real-time anomaly detection applications. It can quickly identify anomalies in network traffic data without the need for extensive computational resources or training time.

Overall, Isolation Forest offers a powerful and efficient approach to anomaly detection in network traffic data. Its ability to handle high-dimensional data, scalability, robustness to noise, and flexibility in data distribution make it a valuable tool for identifying potential security threats and anomalies within network traffic.

**EVALUATION METRIC:**

    **1. Silhouette Score for KMeans:**

The silhouette score measures the quality of clustering by evaluating the separation between clusters and the cohesion within clusters.

A higher silhouette score indicates better-defined clusters, with values closer to 1 indicating dense, well-separated clusters.

A silhouette score of 0.9687 indicates that the clusters formed by KMeans are well-separated and cohesive. This suggests that the algorithm has successfully partitioned the network traffic data into distinct groups based on similarity.

```
Silhouette score: 0.9687799597775262
```

    **2. Silhouette Score for DBSCAN:**

Similar to KMeans, the silhouette score for DBSCAN measures the quality of clustering based on separation and cohesion.

A silhouette score of 0.9449 suggests that DBSCAN has effectively identified dense regions as clusters and separated noise points as outliers. The algorithm has successfully partitioned the network traffic data into clusters with well-defined boundaries.

Despite a slightly lower silhouette score compared to KMeans, DBSCAN still performs well in clustering the network traffic data. The algorithm's ability to adapt to varying cluster densities and identify outliers makes it suitable for detecting anomalies in network traffic instances.

```
Silhouette Score: 0.9449373449177425
```

    **3. Anomaly Scores and Accuracy for Isolation Forest:**

Anomaly scores generated by Isolation Forest represent the degree of abnormality of each data point, with higher scores indicating greater deviation from normal behavior.

The number of instances flagged as anomalies by Isolation Forest is 12597, which matches the expected number based on the contamination rate.

The accuracy of the Isolation Forest model is calculated as the proportion of correctly identified anomalies among all instances. The Isolation Forest model performs well in distinguishing between normal and anomalous network traffic instances. The high accuracy suggests that the algorithm effectively identifies anomalies with a minimal false positive rate, making it suitable for detecting potential security threats in network traffic data.

```
Number of instances flagged as anomalies: 12597
Expected number of anomalies based on contamination rate: 12597
Accuracy of Isolation Forest model: 0.9000023814626944
```

    **4. Anomaly Detection Results for One-Class SVM:**

One-Class SVM classifies instances as either normal or anomalous based on their distance from the learned decision boundary.

The number of instances flagged as anomalies by One-Class SVM is 2569, slightly higher than the expected number based on the nu parameter (2519).

The accuracy of the One-Class SVM model is calculated as the proportion of correctly identified anomalies among all instances. With an accuracy of 0.8980, the model performs well in detecting anomalies in network traffic data.

One-Class SVM also performs well in detecting anomalies in network traffic data. Although the number of instances flagged as anomalies slightly exceeds the expected number based on the nu parameter, the high accuracy suggests that the algorithm effectively distinguishes between normal and anomalous network traffic instances.

```
Number of instances flagged as anomalies: 2569
Expected number of anomalies based on nu parameter: 2519
Accuracy of One-Class SVM model: 0.8980353244691407
```

## COMPARISON BETWEEN THE MODELS:

In the evaluation of anomaly detection models applied to network traffic data, each algorithm demonstrates distinct strengths and characteristics. KMeans and DBSCAN, both clustering algorithms, exhibit high silhouette scores, indicating well-defined clusters within the data. KMeans excels in identifying cohesive clusters, while DBSCAN offers flexibility in handling varying cluster shapes and densities. Isolation Forest, leveraging anomaly scores and accuracy metrics, effectively identifies anomalies with minimal false positives. Its scalability, efficiency, and robustness to noise make it a standout choice for detecting potential security threats. One-Class SVM also performs well, with high accuracy in distinguishing between normal and anomalous instances, albeit slightly exceeding the expected number of anomalies. Each model showcases unique advantages, from KMeans and DBSCAN's clustering capabilities to Isolation Forest's efficiency and One-Class SVM's ability to model complex decision boundaries.

## BEST FITTED MODEL:

Considering all the evaluation metrics and characteristics of each model, Isolation Forest emerges as the best-fitted model for anomaly detection in network traffic data. It demonstrates several key advantages that make it well-suited for this task:

High Accuracy: Isolation Forest achieves a high accuracy of 0.9000, indicating its effectiveness in distinguishing between normal and anomalous network traffic instances with minimal false positives.
Efficiency and Scalability: Isolation Forest is computationally efficient and scalable to large datasets, making it suitable for analyzing extensive network traffic records captured over extended periods.
Robustness to Noise and Outliers: Isolation Forest is robust to noise and outliers in the dataset, effectively isolating anomalies and focusing on the intrinsic structure of the data.
Flexibility in Feature Space: Isolation Forest can operate in high-dimensional feature spaces, capturing complex relationships and interactions between features in network traffic data.
Minimal Parameter Tuning: Isolation Forest requires minimal parameter tuning compared to other algorithms, offering simplicity and ease of use in practice.

Overall, the combination of high accuracy, efficiency, scalability, and robustness to noise positions Isolation Forest as the best-fitted model for anomaly detection in network traffic data. Its ability to effectively detect potential security threats while minimizing false positives makes it a valuable tool for enhancing network security measures.

KMeans achieved a high silhouette score of 0.9687, indicating well-separated and cohesive clusters, which is indeed a significant performance metric. However, when determining the best-fitted model for anomaly detection in network traffic data, it's essential to consider various factors beyond a single evaluation metric. While KMeans performed well in terms of clustering quality, it's crucial to remember that anomaly detection is not solely about clustering. An ideal anomaly detection model should accurately identify anomalies while minimizing false positives and negatives, handle noise and outliers effectively, be scalable to large datasets, and require minimal parameter tuning for practical deployment.

In comparison to KMeans, Isolation Forest offers a comprehensive set of advantages, including high accuracy, efficiency, scalability, robustness to noise and outliers, and minimal parameter tuning. Its anomaly detection capabilities encompass a broader spectrum of requirements for network traffic analysis, making it more suitable as the best-fitted model for anomaly detection in this context.

Therefore, while KMeans' high silhouette score is commendable and indicative of its clustering performance, Isolation Forest's overall performance across multiple metrics makes it the preferred choice for anomaly detection in network traffic data.

In conclusion, the evaluation of anomaly detection models applied to network traffic data highlights the effectiveness of Isolation Forest as the best-fitted model for this task. While KMeans demonstrated a high silhouette score, indicating well-separated clusters, Isolation Forest outperformed other models across multiple evaluation metrics. With its high accuracy, efficiency, scalability, robustness to noise and outliers, and minimal parameter tuning requirements, Isolation Forest excels in accurately identifying anomalies in network traffic data while minimizing false positives. As a result, Isolation Forest emerges as the preferred choice for enhancing network security measures by effectively detecting potential security threats within network traffic.

**CONCLUSION:**

The project aimed to develop and evaluate anomaly detection models for network traffic data, with the ultimate goal of enhancing network security measures. Beginning with data collection from Kaggle, the project utilized various machine learning algorithms, including KMeans, DBSCAN, Isolation Forest, and One-Class SVM, to detect anomalies in network traffic. After preprocessing the data and implementing the models, thorough evaluation was conducted using relevant metrics. The evaluation revealed that while KMeans demonstrated strong clustering performance with a high silhouette score, Isolation Forest emerged as the best-fitted model for anomaly detection. Isolation Forest exhibited high accuracy, efficiency, scalability, and robustness to noise and outliers, making it well-suited for effectively identifying anomalies in network traffic data.

The findings suggest that Isolation Forest offers a comprehensive solution for detecting potential security threats within network traffic, thereby contributing to improved network security measures. Future work could involve further refinement of the models, exploring ensemble methods, or integrating real-time anomaly detection systems into network infrastructures for proactive threat mitigation.

In summary, the project successfully developed and evaluated anomaly detection models for network traffic data, with Isolation Forest identified as the most effective model for enhancing network security by accurately detecting anomalies and potential security threat.

Colab Notebook: [Anomaly Detection in Network Traffic](Anomaly Detection in Network Traffic)