

Anomaly Detection

Provided by: Nour Almulhem

- **Supervised Machine Learning for Anomaly Detection**

Supervised machine learning involves creating a predictive model using a labeled training set that includes both normal and anomalous samples. Common supervised methods include Bayesian networks, k-nearest neighbors, decision trees, supervised neural networks, and support vector machines (SVMs).

The key advantage of supervised models is their potentially higher detection rate compared to unsupervised techniques. This is because they can:

- Provide a confidence score with the model output.
- Incorporate both data and prior knowledge.
- Encode interdependencies between variables.

- **Unsupervised Machine Learning for Anomaly Detection**

Unsupervised methods do not require manually labeled training data. Instead, they operate under the assumption that a small, statistically distinct percentage of network traffic is malicious and abnormal. These methods assume that frequent, similar instances are normal, while infrequent data groups are flagged as malicious.

Popular unsupervised anomaly detection algorithms include Autoencoders, K-means, Gaussian Mixture Models (GMMs), hypothesis test-based analysis, and Principal Component Analysis (PCA).

- **Chosen Algorithms and Their Effectiveness:**

Isolation Forest

Isolation Forests isolate anomalies rather than profiling normal points. The algorithm works by recursively splitting the data and isolating points until separate trees can isolate instances. The IsolationForest classifier from Scikit-Learn can be used in Python. Points requiring fewer splits to isolate are more likely to be anomalies.

IQR Method / Variance-Based

The Interquartile Range (IQR) method is a simple statistical technique that defines anomalies as data points falling below $Q1 - 1.5 \cdot IQR$ or above $Q3 + 1.5 \cdot IQR$, where $Q1$ and $Q3$ are the first and third quartiles. This method can be easily implemented in Python using the Pandas and NumPy libraries.

One-Class SVM

A support vector machine (SVM) can be used for anomaly detection in an unsupervised manner through its one-class extension. One-class SVM learns a decision boundary that envelops most of the regular data points, classifying test points outside this boundary as anomalies. Scikit-Learn provides support for this technique.

Autoencoder

Autoencoders are neural networks that encode and reconstruct input data, trained to minimize reconstruction error. Instances with high reconstruction error at test time are classified as anomalies. Keras and PyTorch libraries in Python support autoencoder-based anomaly detection.

These methods offer a solid starting point for detecting anomalies in Python. The choice of method depends on the use case, data size, and other constraints.