

Machine Learning Anomaly Detection

Azka Javaid

December 12, 2016

Anomaly Detection

- Identify intrusive/anomalous events outside bounds of normal behaviors
- Anomaly Detection System consists of:
 - Data collection
 - Data preprocessing
 - Normal behavior learning phase
 - Anomaly detection
 - Defense response

Machine Learning Overview

- Process of automatic inferring and generalizing a learned model from data
- Branch of larger discipline of artificial intelligence, cognitive computing and deep learning
- Applications include:
 - Building personalized user recommendations
 - Natural language processing (NLP), sentiment analysis, facial and speech recognition for identifying cyber attacks
 - Fraud detection/attack likelihood prediction in dynamic settings (Advanced Persistent Attacks: APT)

Forms of Machine Learning

- Supervised Learning: Data pre-labeled with anomalous and normal features
 - Logistic Regression
 - Decision Trees
 - Artificial Neural Networks (ANN)
- Unsupervised Learning: Pre-existing classification categories not available
 - k-means clustering

HTTP Dataset CSIC 2010

- Developed at the Information Security Institute of CSIC (Spanish Research National Council)
- Address lack of publicly available data to test Firewalls
- Automatic web request traffic to an e-Commerce web application
- Considers static and dynamic anomalous requests:
 - SQL and CRLF injection
 - Buffer overflow
 - Cross-site scripting
- Attacks generated using Paros and Web Application Attack and Audit Framework (W3AF)

HTTP Dataset Feature Description

- HTTP protocol features:
 - Method (GET/POST/PUT)
 - Url
 - Content Length
 - Cookie, Payload (key = value pairs)
- Added features:
 - Index (used to track HTTP packets)
 - Label (anomalous/normal)
 - countPayload (count of total characters in a payload)
 - countJSession (count of unique JSESSIONIDs grouped by url)
 - countIndex (count of unique index values grouped by url)

GOAL: Predict whether an attack is normal or anomalous (label) based on method, contentLength, countPayload, countJSession and countIndex

Findings

- Low (≤ 4.5) countPayload and low (< 225) countJSession indicative of anomalous behavior
- Low (≤ 4.5) countPayload and high (≥ 225) countJSession indicative of normal behavior