Chapter 1

# INTRODUCTION

Computer networks are telecommunication networks which allow computers to exchange data. Network security is implemented on a variety of computer networks to secure daily transactions and communications among businesses, governments and individuals.

## 1.1 Overview

The current network security technologies (such as Firewall, IDS/ IPS) do provide security to networks but have limitations in terms of detecting only the known attack patterns. The system becomes prone to unknown vulnerabilities until the signatures of such attacks are found and configured in the signature databases of network security solutions. The moment attack pattern changes, the signature based solutions fail.

## 1.2 Network Security

Network security[1] consists of the provisions and policies adopted by a network administrator to prevent and monitor unauthorized access, misuse, modification, or denial of a computer network and network-accessible resources. Network security involves the authorization of access to data in a network, which is controlled by the network administrator.

### 1.2.1 Anomaly in Network Security

In network security, **anomaly** (or **outlier** ) is an item, event or observation which does not conform to an expected pattern or other items in a dataset. Typically the anomalous items will translate to some kind of problem such as an attack on the network or a network defect. Anomalies are also referred to as outliers, novelties, noise, deviations and exceptions.

### 1.2.2 Intrusion Detection Systems

An Intrusion Detection System (IDS)[2] is a device or software application that monitors network or system activities for malicious activities or policy violations and produces reports to a management station as illustrated in figure 1.1. Intrusion detection and prevention systems (IDPS) are primarily focused on identifying possible incidents, logging

information about them, and reporting attempts.

In a passive system, the intrusion detection system (IDS) sensor detects a potential security breach, logs the information and signals an alert on the console and/or owner. In a reactive system, also known as an intrusion prevention system (IPS), the IPS auto-responds to the suspicious activity by resetting the connection or by reprogramming the firewall to block network traffic from the suspected malicious source. The term IDPS is commonly used where this can happen automatically or at the command of an operator; systems that both "detect (alert)"and "prevent".
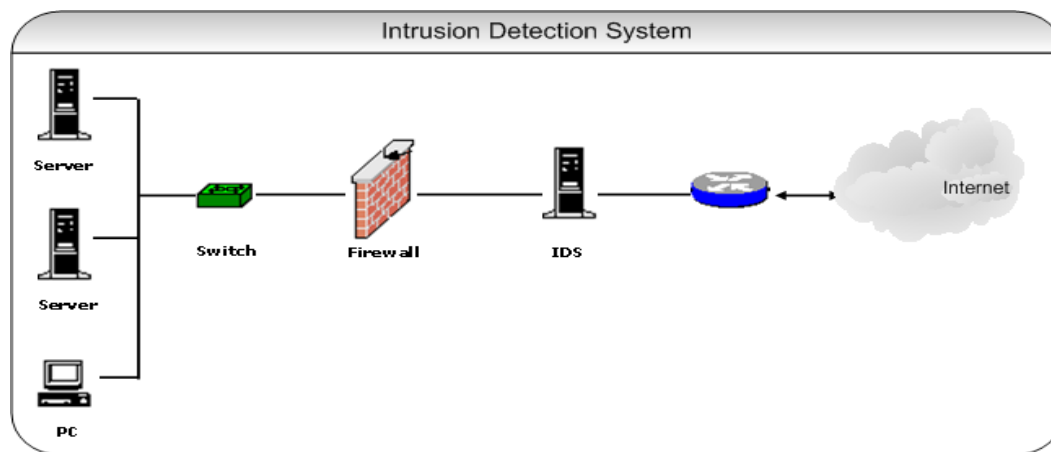


**Figure 1.1 Intrusion Detection Systems**
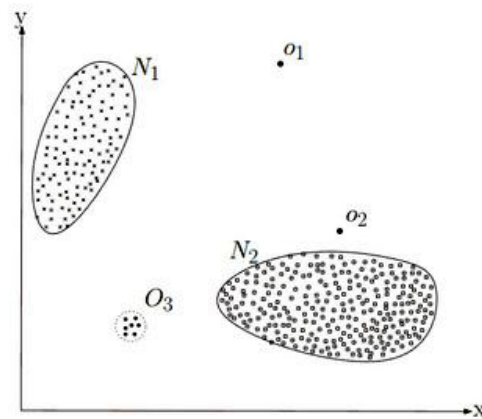
## 1.2.3 Attack Classification

In computer networks, an attack is any attempt to destroy, expose, alter, disable, steal or gain unauthorized access to or make unauthorized use of an asset.[3] Common attack classes and their types are outlined in table 1.1.

| Attack Class | Attack Type |
|---|---|
| Probe | portsweep, ipsweep, quesosatan, msscan, ntinfoscan, lsdomain, illegal-sniffer. |
| DoS | apache2, smurf, neptune, tosnuke, land, pod, back, teardrop, tcpreset, syslogd, crashiis, arppoison, mailbomb, selfping, processtable, udpstorm, warezclient. |
| R2L | dict, netcat, sendmail, imap, ncftp, xlock, xsnoop, sshtrojan, framespoof, ppmacro, guest, netbus, snmpget, ftpwrite, hhtptunnel, phf, named. |
| U2R | sechole, xterm, eject, ps, nukepw, secret, perl, yaga, fdformat, ffbconfig, casesen, ntfsdos, ppmacro, loadmodule, sqlattack. |

**Table 1.1 Attacks present in DARPA 1999 Dataset**

# 1.3 Machine Learning and its Applications

Machine learning, a branch of artificial intelligence, concerns the construction and study of systems that can learn from data. For example, a machine learning system could be trained on network traffic to learn to distinguish between anomalous and normal traffic. After learning, it can used to classify the incoming traffic into anomalous and normal traffic as shown in figure 1.2.



A simple example of anomalies in a 2-dimensional data set.

**Figure 1.2 Example of Anomalies**

Machine Learning is used in a variety of fields, a few of which are mentioned below.

   i. **Network Anomaly Detection:** Statistical machine learning is one class of statistical approaches used in the field of Network Anomaly Detection. A choice of a learning scenario depends on the information available in measurements. For example, a frequent scenario is that there are only raw network measurements available and thus unsupervised learning methods are used. If additional information is available, e.g. from network operators, known anomalies or normal network behaviors, learning can be done with supervision. A choice of a mapping depends on the availability of a model, the amount and the type of measurements, and complexity of learning algorithms[4].

   ii. **Evaluating Learned Knowledge:** Rules induced from training data are not necessarily of high quality. The performance of knowledge acquired in this

way is an empirical question that must be answered before that knowledge can be used on a regular basis. One standard approach to evaluation involves dividing the data into two sets, training on the first set, and testing the induced knowledge on the second. One can repeat this process a number of times with different splits, and then average the results to estimate the rules' performance on completely new problems. Kibler and Langley (1988) experimental methods of this sort for a broad class of learning algorithms. An important part of the evaluation process is experts' examination of the learned knowledge. Evans and Fisher (1994) encourage an iterative process in developing a fielded application.[5]

## 1.4 Types of Intrusion Detection Systems

An Intrusion can be termed as any activity which is aimed at disrupting or denying a service by gaining unauthorized access. Various attacks like DDoS, Identity spoofing, Ping of Death, Buffer Overflow etc. are used to intrude and disrupt a network. Intrusion Detection Systems are used to prevent such intrusions by detecting them at an earlier stage.

All Intrusion Detection Systems use one of two detection techniques:

i.  Signature Based Detection

ii. Statistical Anomaly Based Detection

### 1.4.1 Signature Based Detection

Signature-based IDS monitors packets in the network, and compares them with pre-configured and pre-determined attack patterns, known as signatures.[6]

### 1.4.2 Statistical Anomaly Based Detection

Anomaly based detection techniques model a "normal" scenario and any deviation from that is considered an anomaly. Anomaly based detection is preferred over Signature based detection due to their ability to detect novelty attacks.[6]

## 1.5 Machine Learning

Machine learning is a branch of artificial intelligence, concerned with the construction and study of systems that can learn from data. Machine learning focuses on the

development of computer programs that can teach themselves to grow and change when exposed to new data as illustrated in figure 1.3.[7]
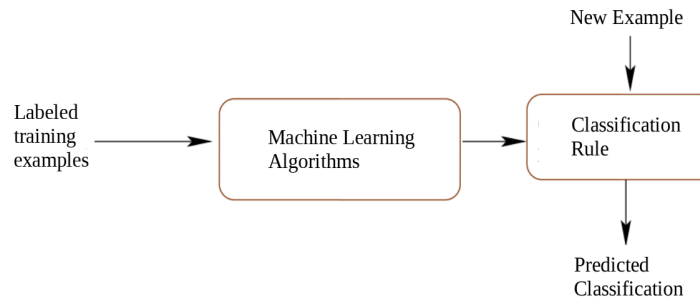


**Figure 1.3 Machine Learning System**

Machine learning algorithms are classified into 3 types-

    i.   Supervised Learning

    ii.   Unsupervised Learning

    iii.  Reinforcement Learning

## 1.5.1 Supervised Learning

In supervised learning, the model defines the effect one set of observations, called inputs, has on another set of observations, called outputs.[7] As shown in Figure 1.4, the inputs are assumed to be at the beginning and outputs at the end of the causal chain. The models can include mediating variables between the inputs and outputs. A supervised learning algorithm makes it possible, given a large data set of input/output pairs, to predict the output value for some new input value that you may not have seen before.
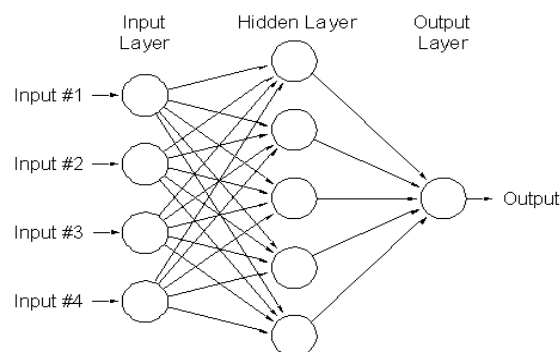


**Figure 1.4 Supervised Learning Model**

## 1.5.2 Unsupervised Learning

In unsupervised learning, all the observations are assumed to be caused by latent variables, that is, the observations are assumed to be at the end of the causal chain.[7] With unsupervised learning it is possible to learn larger and more complex models than with supervised learning. The learning can proceed hierarchically from the observations into ever more abstract levels of representation as shown in figure 1.5. Each additional hierarchy needs to learn only one step and therefore the learning time increases (approximately) linearly in the number of levels in the model hierarchy unlike supervised learning where the learning task increases exponentially. Unsupervised learning can be used for bridging the causal gap between input and output observations. The latent variables in the higher levels of abstraction are the causes for both sets of observations and mediate the dependence between inputs and outputs.
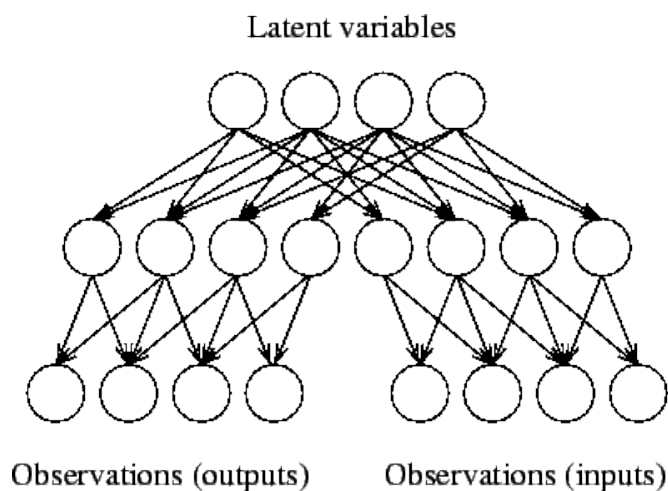


**Figure 1.5 Unsupervised Learning Model**

## 1.5.3 Reinforcement Learning

Reinforcement learning is learning what to do - how to map situations to actions -  as to maximize a numerical reward signal. A Reinforcement learning agent learns from the consequences of its actions, rather than from being explicitly taught and it selects its actions on basis of its experiences (exploitation) and also by new choices (exploration), which is essentially trial and error learning as illustrated in figure 1.6.
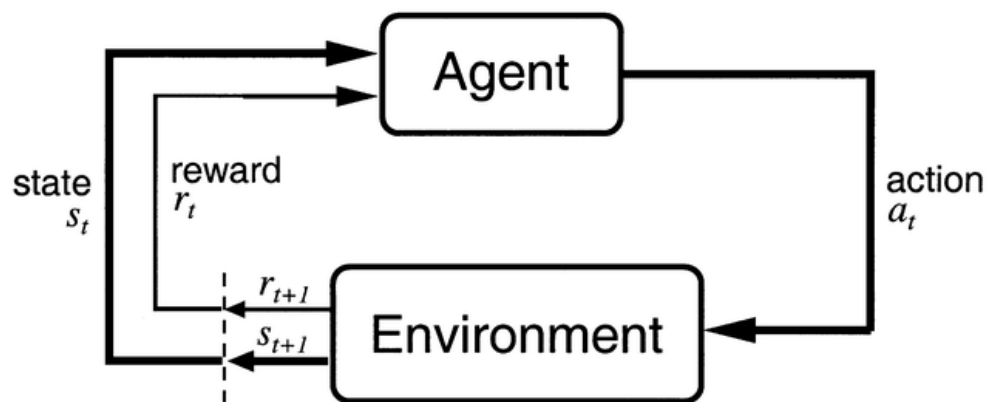
**Figure 1.6 Reinforcement Learning Model**

# 1.6 Types of Anomaly Based Detection

We now review the different techniques of Anomaly based IDS. They are classified as shown in figure 1.7. The most important ones are-

i.   Statistical Anomaly Detection

ii.  Data Mining Based Detection
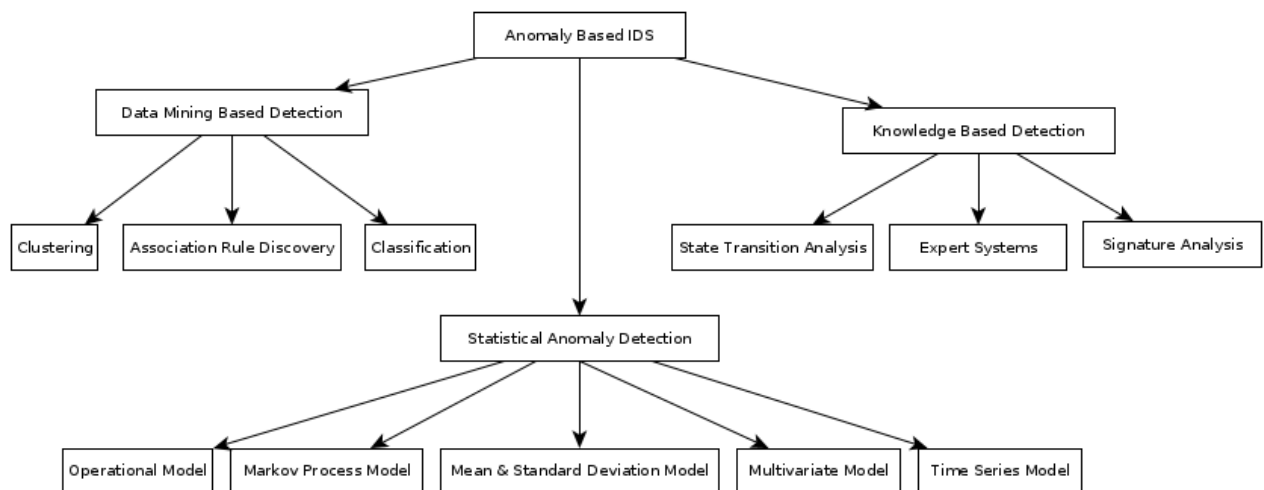
iii. Knowledge Based Detection



**Figure 1.7 Taxonomy of Anomaly Based IDS**

## 1.6.1 Statistical Anomaly Detection

Statistical based anomaly detection techniques use statistical properties and statistical tests to determine whether "Observed behavior" deviate significantly from the "expected

---

behavior". Statistical based anomaly detection techniques use statistical properties (e.g., mean and variance) of normal activities to build a statistical based normal profile and employ statistical tests to determine whether observed activities deviate significantly from the normal profile. The IDS goes on assigning a score to an anomalous activity. As soon as this score becomes greater certain threshold, it will generate an alarm.[8] Statistical Anomaly Based IDS can further be classified into the following categories-

- **Operational Model:** This model is based on the operational assumption that an anomaly can be identified through a comparison of an observation with a predefined limit. Based on the cardinality of observation that happens over a time an alarm is raised. The operational model is most applicable to metrics where experience has shown that certain values are frequently linked with intrusions. For example, an event counter for the number of password failures during a brief period, where more than say 10, suggest a failed log-in.

- **Markov Process Model:** The Markovian Model is used with the event counter metric to determine the normalcy of a particular event, based on the events which preceded it.[8] The model characterizes each observation as a specific state and utilizes a state transition matrix to determine if the probability of the event is high (normal) based on the preceding events. This model is basically used in two main approaches: Markov chains and Hidden Markov models. A Markov chain keeps track of an intrusion by examining the system at fixed intervals and maintains the record of its state. If the state change takes place it computes the probability for that state at a given time interval. If this probability is low at that time interval then that event is considered as an anomalous.

- **Mean and Standard Deviation Model:** The Mean and Standard Deviation Model is based on the traditional statistical determination of the normalcy of an observation based on its position relative to a specified confidence range.[8] In this sub-model, event that falls outside the set interval will be declared as an anomalous.

- **Multivariate Model:** This model can be applied to intrusion detection for monitoring and detecting anomalies of a process in an information system.[8] This model is

similar to the mean and standard deviation model except that it is based on correlations among two or more features. This model would be useful in the situation where two or more features are related. This model permits the identification of potential anomalies where the complexity of the situation requires the comparison of multiple parameters.

- **Time Series Model:** The Time Series Model, attempts to identify anomalies by reviewing the order and time interval of activities on the network.[8] If the probability of the occurrence of an observation is low, then the event is labeled as abnormal. This model provides the ability to evolve over time based on the activities of the users. Anomalies in time series data are data points that significantly deviate from the normal pattern of the data sequence.

## 1.6.2 Data Mining Based Approach

As IDS can only detect known attacks, but it cannot detect insider attacks, the better solution for an IDS can be Data Mining at its core and is defined as "the process of extracting useful and previously unnoticed models or patterns from large data stores".[8] The data-mining process tends to reduce the amount of data that must be retained for historical comparisons of network activity, creating data that is more meaningful to anomaly detection.

Data Mining based approach is classified into following techniques:

- **Clustering:** Clustering[8] is an unsupervised technique for finding patterns in unlabeled data with many dimensions (number of attributes). Mostly k-means clustering is used to find natural groupings of similar instances. Records that are far from any of these clusters indicate unusual activity that may be part of a new attack.

- **Association Rule Discovery:** Association rules mining finds correlation between the attributes. This technique was initially applied to the so-called market basket analysis, which aims at finding regularities in shopping behavior of customers of supermarkets. This method is usually very slow and its being replaced by other powerful techniques like clustering and classification.[8]

- **Classification:** Classification[8] is one of the main techniques used in data mining. Its main goal is to learn from class-labeled training instances for predicting classes of new or previously unseen data. Instances in the training set have labels. New data is classified based on the training set. Basically a classification tree (also called as decision tree) is built to predict the category to which a particular instance belongs.

## 1.6.3 Knowledge Based Detection Techniques

Knowledge based detection Technique[8] can be used for both signature based IDS as well as anomaly based IDS. It accumulates the knowledge about specific attacks and system vulnerabilities. It uses this knowledge to exploit the attacks and vulnerabilities to generate the alarm. Any other event that is not recognized as an attack is accepted. Therefore the accuracy of knowledge based intrusion detection systems is considered good. Knowledge based detection Technique can further be classified as:

i. **State Transition Analysis:** This technique describes the attacks with a set of goals and transitions, and represents them as state transition diagrams.[8] State transition diagram is a graphical representation of the actions performed by an intruder to archive a system compromise. In state transition analysis, an intrusion is viewed as a sequence of actions performed by an intruder that leads from some initial state on a computer system to a target compromised state. State transition analysis diagrams identify the requirements and the compromise of the penetration. They also list the key actions that have to occur for the successful completion of an intrusion.

ii. **Expert Systems:** The expert system[8] contains a set of rules that describe attacks. Audit events are then translated into facts carrying their semantic significance in the expert system, and the inference engine draws conclusions using these rules and facts. This method increases the abstraction level of the audit data by attaching semantic to it. It also encodes knowledge about past intrusions, known system vulnerabilities and the security policy. As information is gathered, the expert system determines whether any rules have been satisfied.

iii. **Signature Analysis:** Signature analysis follows exactly the same knowledge-acquisition approach as expert systems, but the knowledge acquired is exploited in a different way.[8] The semantic description of the attacks is transformed into information that can be found in the audit trail in a straightforward way. For example, attack scenarios might be translated into the sequences of audit events they generate, or into patterns of data that can be sought in the audit trail generated by the system.