

Network Traffic Classification for Cybersecurity and Monitoring

Introduction

Network traffic classification for cybersecurity and monitoring is an essential process that involves analyzing and categorizing network packets to identify various types of traffic. This classification helps in monitoring network behavior, detecting anomalies, enhancing security measures, and optimizing network performance.

String searching : Aho-corasick

String searching is a fundamental problem in computer science where the goal is to find one or more occurrences of a substring (pattern) within a larger string (text).

nDPI implementation : String matching in nDPI is employed in the following contexts:

- **Protocol Identification**: Uses string matching to detect and classify network protocols.
- **Payload Analysis**: Scans packet payloads for specific byte sequences.
- **Application Detection**: Identifies applications by matching known signatures.
- **Intrusion Detection**: Detects malicious activity by comparing payloads to known attack signatures.

Aho-corasick : The Aho-Corasick algorithm is designed for multiple pattern searching. It constructs a finite state machine (trie) with failure links to handle pattern mismatches efficiently.

- **Complexity**: $O(n+m+z)$, where z is the total number of matches.
- **Usage**: Suitable for applications like network intrusion detection and text processing where multiple patterns need to be searched simultaneously.

IP Matching : Radix tree

IP matching involves checking if an IP address belongs to a specific subnet or range. This is fundamental in tasks such as routing, firewall filtering, and IP address management.

A radix tree is a data structure that is commonly used to store and efficiently look up IP addresses. It is a compressed version of a binary trie, where each node that is the only child is merged with its parent.

Structure

- **Nodes:** Represent bits of the IP address.
- **Edges:** Transition from one bit to the next.
- **Leaf Nodes:** Store IP address prefixes.

Operations

- **Insert:** Add an IP address or prefix to the tree.
- **Lookup:** Find the longest matching prefix for a given IP address.

Probabilistic Counting: HyperLogLog

HyperLogLog is a probabilistic data structure used to estimate the cardinality of a set. It improves probabilistic counting by hashing every element, and counting the amount of 0s to the left of such hash.

How HyperLogLog Works ?

1. **Hashing:** Each element in the data stream is hashed to a uniformly distributed random value. For a 32-bit hash, this would produce a number in the range $[0, 2^{32}-1]$.
2. **Splitting Hash Values:**
 - The hash value is divided into two parts: the first part is used to determine the register index, and the second part is used to count leading zeros.
 - For example, with 2^b registers, the first b bits of the hash value determine the register, and the remaining bits are used to count leading zeros.
3. **Updating Registers:**
 - For each element, compute the hash and split it.
 - Determine the register index using the first part of the hash.
 - Update the register with the maximum number of leading zeros observed in the second part.
4. **Estimate Calculation:**
 - The distinct count estimate is derived from the harmonic mean of the values in the registers.
 - Bias correction and scaling factors are applied to refine the estimate.

Anomaly Detection :

Anomaly detection, the process of identifying data points that deviate significantly from the norm, is a critical task in various domains such as cybersecurity, finance, healthcare, and industrial monitoring.

Outliers are data points that significantly deviate from the majority of a dataset. These anomalies can be higher or lower than the other values and often indicate errors, variability in the data, or novel insights

nDPI implements three “smoothing” functions for data forecast :

- **Single Exponential Smoothing:** For data without trend or seasonality.
- **Double Exponential Smoothing:** For data with a trend.
- **Triple Exponential Smoothing (Holt-Winters):** For data with trend and seasonality.

Data Comparison: Binning

Binning is a technique used in data analysis and preprocessing to group continuous data into discrete intervals or "bins." This approach can be useful for various purposes, including data summarization, visualization, and preparation for algorithms that require categorical input.

Purpose of Binning

- **Data Reduction:** Simplifies data by grouping it into fewer categories.
- **Data Visualization:** Makes patterns more apparent by reducing the noise of raw data.
- **Preprocessing:** Transforms continuous data into a categorical format suitable for certain algorithms.

Types of Binning

- **Equal-Width Binning:** Divides the data range into intervals of equal size.
 - **Example:** For data ranging from 0 to 100, creating 10 bins each of width 10 (0-10, 10-20, etc.).
- **Equal-Frequency Binning:** Divides the data so that each bin contains approximately the same number of data points.
 - **Example:** Sorting the data and then creating bins so that each bin has an equal number of data points.
- **Custom Binning:** Defines bins based on specific criteria or domain knowledge.
 - **Example:** Age groups like 0-18, 19-35, 36-55, 56+.
- **Clustering-Based Binning:** Uses clustering algorithms to create bins based on the distribution of the data.
 - **Example:** K-means clustering to group data into clusters, which can then be used as bins.

Applications of Binning

- **Histogram Creation:** To visualize the distribution of data.

- **Feature Engineering:** To prepare data for machine learning algorithms.
- **Data Aggregation:** To summarize and report data in a more meaningful way.