

BigFlow: Real-time and Reliable Anomaly based Intrusion Detection for High-Speed Networks

Future Generation Computer Systems

2019

73 citations

Keywords: datastreaming

روشهای فعلی تشخیص برای داده‌های حجیم، بسته‌ها را ضبط کرده و آنها را به یک فایل سیستم دیگر که از نوع HDFS برای آنالیز می‌فرستند. الگوریتم‌های یادگیری ماشین فعلی در سرعت‌های بالا و عملکرد بلادرنگ^۱ مشکل دارند روشهای یادگیری ماشین غیرنظارتی^۲ برای کشف آنومالی‌ها، به دلیل ضبط بسته‌ها در یک بازه زمانی و سپس تشخیص بی‌نظمی نیز مشکل حافظه خواهند داشت و همچنین نرخ مثبت کاذب بالایی خواهند داشت. لذا روشهای نظارتی پیشنهاد شدند، که بر اساس یک مدل به دست‌آمده (بر اساس یک پروسه تمرین^۳ اغلب طولانی و هزینه‌بر) ترافیک‌های بد و نابد را با استفاده از یک الگوریتم classifier دسته‌بندی می‌کنند. اما با توجه به ماهیت متغیر ترافیک، این مدل را باید مرتباً بازطراحی کرد. در این آزمایش بررسی می‌کنیم که دقت این سامانه‌ها در طول سال به ۲۳ درصد کاهش می‌یابد. لذا برای نشان دادن آن به عنوان یک سیستم قابل اتکا باید مدام آموزش و تست شود.

راهکار ما در این مقاله از این روش استفاده می‌کند که آیا خروجی classifier که همان تعیین مهاجم بودن یا سالم بودن ترافیک است، باید قبول شود (rejecting low confidence classifications)؟) با کمک event class(normal or attack) (probability(confidence)). مثلاً اگر رویدادی مهاجم طبقه‌بندی شده بود، تنها در صورتی می‌تواند قبول شود که معیار اطمینان آن بالای ۹۰٪ باشد. اگر یکی قبول نشد، ادمین می‌فهمد که رفتار شبکه تغییر پیدا کرده است. تغییر در رفتار شبکه می‌تواند به خاطر ظهور سرویس‌های جدید باشد. همچنین راهکارمان یک stream learning می‌باشد که ترافیک را در لحظه بررسی خواهد کرد و به روزرسانی افزایشی مدل بر اساس نمونه‌های رد شده را خواهیم داشت. این مدل‌های افزایشی باعث افزایش دقت و کاهش زمان یادگیری (به دلیل دور نیانداختن مدل قبلی) در شبکه‌های پهن باند خواهند شد. اما مشکلشان این است که نظارتی هستند یعنی نیاز به داده‌های از قبل طبقه‌بندی شده^۴ خواهند داشت. روش ما می‌خواهد مداخله انسان را کاهش دهد و همچنین میزان دپتایی که باید ضبط شود را کم کند.

نوآوری‌های ما در این مقاله:

¹ Real-time

² Unsupervised

³ training

⁴ Classified

- تولید اولین دیتاست برای بررسی IDS در مدت زمان طولانی یکسال که شامل رکوردهای برچسب گذاری شده (نرمال یا مهاجم) می باشد که هر بسته ۱۵۸ تا feature خواهد داشت (براساس دیتاست MAWIFlow می باشد). دیتاست های که برای بررسی سامانه های تشخیص نفوذ استفاده می شوند باید دارای ویژگی های زیر باشند: validity, realism. با این که پکت ها پیلودشان حذف شده است و داده های حساس موجود در سرایندها حذف شده اند اما همچنان بازسازی جریان ها ممکن است چون که feature ها دستکاری نشده اند، prior labeling, high variability, reproducibility and public availability
- تست الگوریتم های classifier مختلف با این دیتاست (الگوریتم های decision tree, random forest, gradient boosting, hybrid). برای هر کدام نیز دو گونه آپدیت (آموزش مدل) مدل در نظر گرفتیم: بدون آپدیت و آپدیت هفتگی. از Apache Spark برای پیاده سازی و ارزیابی این الگوریتم ها استفاده کرده ایم.
- طراحی bigflow که آن را در یک بازه زمان یکساله بررسی کرده و دقت بالایی و میزان استفاده کم از منابع و همچنین زمان کم برای آموزش دارا می باشد.

هنگام استفاده از روش های یادگیری ماشین، رفتار ترافیک بر اساس یک سری ویژگی هایی که از آن استخراج می شود، نمایش داده می شود:

Table 1

Network-level feature set used in the experiments throughout this work [18].

Type	Grouping	Features
Host-based	Host to All	Number of Packets, Number of Bytes, Average Packet Size, Percentage of Packets (PSH Flag), Percentage of Packets (SYN and FIN Flags), Percentage of Packets (FIN Flag), Percentage of Packets (SYN Flag), Percentage of Packets (ACK Flag), Percentage of Packets (RST Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Time Exceeded Flag), Percentage of Packets (ICMP Unreachable Flag), Percentage of Packets (ICMP Other Types Flag), Average Packet Size, Throughput in Bytes, Protocol
Flow-based	Source to Destination	Number of Packets, Number of Bytes, Average Packet Size, Percentage of Packets (PSH Flag), Percentage of Packets (SYN and FIN Flags), Percentage of Packets (FIN Flag), Percentage of Packets (SYN Flag), Percentage of Packets (ACK Flag), Percentage of Packets (RST Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Time Exceeded Flag), Percentage of Packets (ICMP Unreachable Flag), Percentage of Packets (ICMP Other Types Flag), Throughput in Bytes
	Destination to Source	Number of Packets, Number of Bytes, Average Packet Size, Percentage of Packets (PSH Flag), Percentage of Packets (SYN and FIN Flags), Percentage of Packets (FIN Flag), Percentage of Packets (SYN Flag), Percentage of Packets (ACK Flag), Percentage of Packets (RST Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Time Exceeded Flag), Percentage of Packets (ICMP Unreachable Flag), Percentage of Packets (ICMP Other Types Flag), Throughput in Bytes
	Both	Number of Packets, Number of Bytes, Average Packet Size, Percentage of Packets (PSH Flag), Percentage of Packets (SYN and FIN Flags), Percentage of Packets (FIN Flag), Percentage of Packets (SYN Flag), Percentage of Packets (ACK Flag), Percentage of Packets (RST Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Redirect Flag), Percentage of Packets (ICMP Time Exceeded Flag), Percentage of Packets (ICMP Unreachable Flag), Percentage of Packets (ICMP Other Types Flag), Throughput in Bytes

BigFlow بر پایه پلتفرم های پردازش استریمی^۵ بنیان شده است که اینها اطلاعاتی را که دریافت میکنند توسط یک سری PE (Processing Element) پردازش می کنند. ارتباطات و نحوه پیغام دهی بین PE ها سه حالت دارد: shuffle : که با احتمال یکسان به یک PE دیگر می فرستد، keyed: که پیغام ها را مثلن بر اساس آدرس IP گروه بندی کرد و آن را به PE مربوطه ارسال می کند، broadcast: که به تمام PE های همونوع ارسال می کند. بلادرنگی زمانی برقرار می شود که هر PE حافظه محدود داشته باشد و به صورت موازی چندین مورد از آنها باهم پردازش داشته باشند

روش کلی کار BigFlow به دو قسمت تقسیم می شود: (تکته! event (رویداد): یکای چیزی که باید آنالیز شود و می تواند از نوع packet header یا netflow record باشد)

⁵ Apache Flink & Apache Storm

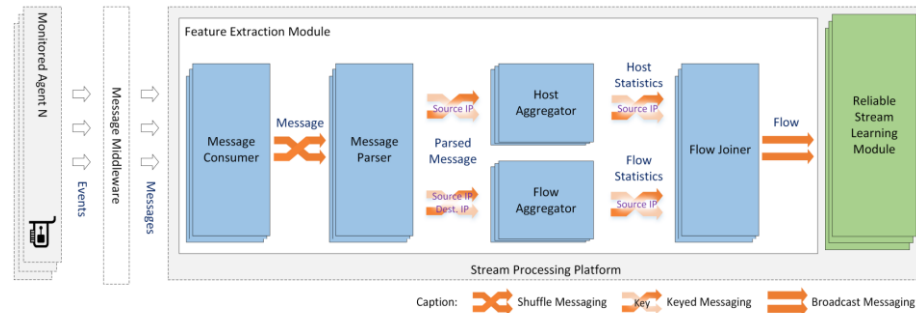


Figure 2 – BigFlow real-time feature extraction module architecture for high-speed networks.

Feature extraction: یک فریمورک پردازش استریم می باشد. آمار جریان به صورت بلادرنگ که **feature vector** نامیده می شود (در یادگیری ماشین **event** یا **instance** گویند) را به دست می آورد و شامل اطلاعات ترافیک ردوبدل شده بین دو هاست در یک بازه زمانی^۶ می باشد. روش های قدیمی استخراج اطلاعات برای **classifieier** در پهن باند کاربرد ندارند. Bigflow ۱۵۸ ویژگی که در سطح میزبان (مثلاً اطلاعات دریافتی و یا ارسالی یک هاست مشخص) و جریان (بین دو هاست مثلاً میانگین حجم بسته های تبادل شده بین دو هاست) می باشند، استخراج می کند. این بخش از اجزای زیر تشکیل شده است:

Monitored Agent: هاست ها، سوئیچ ها، روترها. مثلاً سوئیچ ها هدرهای بسته ای ارسال می کنند و روترها **NetFlow** record ارسال می کنند. رویدادها را از طریق **Message Middleware** به فریمورک ارسال می کنند.

MessageMiddleWare: همه رویدادها از طریق اینترفیس این می روند.

MessageConsumer: دریافت تمامی انواع رویدادهای موجود و سپس ارسال آنها به صورت **shuffle** توسط **PE** های **Stream Processing** به ماژول اصلی.

MessageParser: بر اساس نوع رویداد دریافتی یکسری اطلاعات از آن استخراج می کند: **event** (**source/field/type(netflow record or network packet)** روی **PE** های مختلفی بالانس شده است).

Aggregator: از دو بخش اصلی تشکیل شده است. کار اصلی استخراج ویژگی را این دو انجام میدهند. هر دو از طریق **keyed** پیغام ها رو دریافت می کنند. بر اساس کلیدی که به دست آمده ، به **PE** مشخصی ارسال می کنند (هر کدام از این بخشها از یکسری **PE** تشکیل شده اند):

- **HostAggregator:** کلید آن از هش آدرس آی پی مبدا به دست می آید
- **FlowAggregator:** کلید آن از طریق **XOR** بر روی آدرس مبدا و مقصد به دست می آید.

هر کدام از این ها یک قسمتی از مقادیر هش را مسیول می باشند.

BigFlow برای این که مقادیر رویدادهای مربوط به یک گروه را محاسبه کند، آنها را به یک بازه زمینی تقسیم می کند که ماژولهای **Tumbling Window** نامیده می شوند. هر کدام از این ماژول ها مقادیر ویژگی ها را برای یک بازه زمانی خاصی و رویداد دریافتی

⁶ Time interval

آپدیت یا ذخیره می‌کنند زمانی که **Tumbling Window** منقضی شد، مقادیر ویژگی‌ها با فرمت آمار هاست یا جریان به ماژول بعدی ارسال می‌کند. خاصیت این ماژول این است که ابتدا مطمئن می‌شود که تمامی جریان‌ها منقضی می‌شوند

FlowJoiner: این ماژولها تمامی ویژگیهای انواع مختلف رویدادها را دریافت کرده و آنها را در استریم واحد جریان پیوند می‌کنند. باید توجه داشت که یک هاست واحد ممکن است شامل چندین ویژگی آماری جریانی در یک **Tumbling Window** داشته باشد در حالی که تنها یک ویژگی آماری هاستی داشته باشد. (مثلاً هاستی که از سرویس‌های چندین هاست دیگر استفاده می‌کند). پس وظیفه این ماژول این است که همه اینها را ذخیره کند و در آینده به هم متصل کند.

Reliable stream learning: بر اساس ورودی **feature vector** از ماژول قبلی آن را به مهاجم یا طبیعی **classify** می‌کند. **Stream learning classifier** به همراه یک ماژول **verifier** می‌باشد. زمانی که خروجی **classifier** تایید نشد، اون رویداد را ذخیره می‌کند تا بعداً توسط یک عامل انسانی (که بر اساس اطلاعاتی که از رفتار جدید شبکه به دست می‌آورد مثلاً سرویس‌های جدید و یا سایت‌های CVE) برچسب گذاری شود. سپس از این نمونه برچسب گذاری شده برای به روزرسانی افزایشی **classifier** استفاده می‌شود.

این ماژول در ابتدا با یک دیتاست آموزش داده می‌شود و **classifier model** ساخته می‌شود و مقدار آستانه برای هر کلاس نرمال یا مهاجم تعیین می‌شود.

پیاده سازی: کل طرح ما روی **Apache Flink** که فریمورک مخصوص پردازش استریم می‌باشد پیاده سازی شده است و همچنین مکانیزم **Tumbling Window** نیز در آن شبیه‌سازی شده است و زمان آن نیز به حالت بهینه ۱۵ ثانیه تنظیم شده است. پیام رسانی **keyed** نیز با استفاده از اینترفیس **KeySelector** آن پیاده سازی شده است. (برای آشنایی با بخش‌های دیگر به مقاله مراجعه شود)

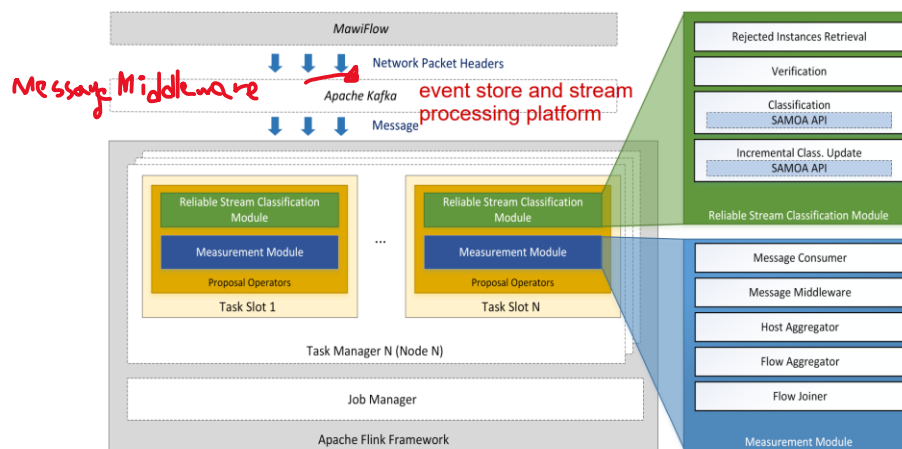


Figure 5 – BigFlow architecture.

ارزیابی: در ابتدا ماژول **Reliable Stream Classification** با معیار دقت در طول زمان، با دیتاست **MAWIFlow** ارزیابی شده است و سپس کارایی و مقیاس پذیری **BigFlow** و همچنین هزینه به روزرسانی ماژول **stream learning** بررسی شده است.

راهکارهای پیشنهادی برای بهبود:

استفاده از اسلج‌ها برای ذخیره ویژگی‌ها؛

تعیین خودکار اینکه که یک event مهاجم است یا نه؛