

Universidad Internacional de la Rioja (UNIR)

Escuela Superior de Ingeniería y Tecnología

Máster Universitario en Análisis y Visualización de Datos Masivos

Caracterización de equipos informáticos según su comportamiento en una red empresarial

Trabajo Fin de Máster

Presentado por: Javier Artiga Garijo

Director: Luis Miguel Garay Gallastegui

Ciudad: Logroño

Fecha: [Fecha de Entrega]

Índice general

1.	Introducción	1
	1.1. Motivación	1
	1.2. Objetivos	2
	1.3. Estructura del documento	2
2.	Objetivos y Metodología	4
	2.1. Objetivo General	4
	2.2. Objetivos Específicos	4
	2.3. Metodología de Trabajo	5
3.	Estado del Arte	7
	3.1. Aprendizaje automático en la clasificación de tráfico	7
	3.2. Detección de anomalías sobre actividad de red	10
	3.3. Clustering	12
4.	Desarrollo	13
	4.1. Presentación del entorno	13
	4.2. Extracción y filtrado	13
	4.3. Preprocesado para el clustering	19
	4.4. Análisis de datos	19
	4.5. Selección de características	19
	4.6. Parametrización	19
5.	Resultados	20
6.	Conclusiones	21
Bi	bliography	22

	Máster	Universitario	en Análisis v	/ Visualización	de	Datos	Masivo
--	--------	---------------	---------------	-----------------	----	--------------	--------

av	ier	Artig	ia	Gari	ic

Appendices	24
Apéndice A. Fusce viverra lectus	25

Resumen

Lorem ipsum dolor sit amet, consectetur adipiscing elit, sed do eiusmod tempor incididunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrud exercitation ullamco laboris nisi ut aliquip ex ea commodo consequat. Duis aute irure dolor in reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint occaecat cupidatat non proident, sunt in culpa qui officia deserunt mollit anim id est laborum.

Palabras clave: Plabra clave 1, Palabra clave 2, Palaba Clave N

Abstract

Proin tincidunt enim nec fringilla euismod. Quisque id efficitur sapien. Sed hendrerit, nisl id efficitur elementum, arcu ex efficitur ipsum, sed auctor metus nibh id leo. Maecenas eu sem tortor. Etiam accumsan bibendum ante vitae auctor. Donec eget dolor gravida, tincidunt diam non, ornare orci. Pellentesque ornare blandit eros, sed maximus neque varius non.

Keywords: Plabra clave 1, Palabra clave 2, Palabra Clave N

Índice de figuras

Índice de Tablas

Introducción

[Los fragmentos entre corchetes son ideas que ampliaré próximamente]

[El primer capítulo es siempre una introducción. En ella debes resumir de forma esquemática pero suficientemente clara lo esencial de cada una de las partes del trabajo. La lectura de este primer capítulo ha de dar una primera idea clara de lo que se pretendía, las conclusiones a las que se ha llegado y del procedimiento seguido.]

[Como tal, es uno de los capítulos más importantes de la memoria. Las ideas principales a transmitir son la identificación del problema a tratar, la justificación de su importancia, los objetivos generales (a grandes rasgos) y un adelanto de la contribución que esperas hacer. Típicamente una introducción tiene la siguiente estructura:

- Motivación: ¿Cuál es el problema que quieres tratar? ¿Cuáles crees que son las causas?
 ¿Por qué es relevante el problema?
- 2. Planteamiento del trabajo (sección 1.2): ¿Cómo se podría solucionar el problema? ¿Qué es lo que se propone? Aquí describes tus objetivos en términos generales ("mejorar el aprendizaje de idiomas")
- 3. Estructura del trabajo (sección 1.3): Aquí describes brevemente lo que vas a contar en cada uno de los capítulos siguientes.

1.1. Motivación

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Donec tincidunt libero viverra mi elementum pellentesque. Mauris sollicitudin elementum turpis. Nunc dictum lectus nec ligula fringilla auctor. Proin ut sem felis. Sed suscipit maximus lorem, nec vulputate nisi. Donec eu

interdum lectus 1. Etiam viverra nec nulla vel facilisis. Vivamus purus lacus, laoreet id justo eu, euismod commodo justo.

1.2. Objetivos

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Donec tincidunt libero viverra mi elementum pellentesque. Mauris sollicitudin elementum turpis. Nunc dictum lectus nec ligula fringilla auctor. Proin ut sem felis. Sed suscipit maximus lorem, nec vulputate nisi. Donec eu interdum lectus 1. Etiam viverra nec nulla vel facilisis. Vivamus purus lacus, laoreet id justo eu, euismod commodo justo.

Vestibulum risus eros, fringilla quis tristique quis, tincidunt et dui. Sed et magna blandit, sagittis nibh dictum, eleifend magna. Suspendisse a nulla a augue ultrices molestie sit amet et magna. Vestibulum molestie metus id lorem bibendum rhoncus. Maecenas sit amet massa pretium, commodo nunc id, viverra augue. Curabitur nec ultricies sem. Donec congue lectus lorem, vel laoreet arcu mollis volutpat. Duis sem metus, consectetur ac diam et, pretium congue tortor. Quisque elementum mollis enim, nec blandit mauris feugiat quis.

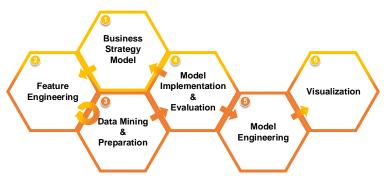


Figura 1: Duis tristique velit velit. Curabitur auctor, nibh.

1.3. Estructura del documento

Lorem ipsum dolor sit amet (entry2013one), consectetur adipiscing elit. Phasellus eu risus convallis, cursus ante quis, condimentum nibh. Cras vulputate suscipit tortor eu pulvinar. Aenean aliquet felis lacinia aliquet interdum. In at diam sed lacus porta efficitur. Nunc hendrerit, felis id venenatis fringilla, augue justo egestas ipsum, vel pretium est lorem at arcu. Nullam finibus mauris et maximus viverra. Etiam sagittis ligula vel aliquam vehicula. Aenean consequat condimentum enim, eget posuere lectus dapibus sed. Suspendisse mollis entry2002two (entry2002two) eu orci eu tempor. Aliquam tempus enim eget egestas congue. Morbi nisi metus, suscipit in scelerisque quis, accumsan sit amet nulla. Sed id imperdiet quam, eget

commodo sem. Donec et molestie libero. Donec et quam mi. Vivamus consectetur est turpis, et ultricies leo dictum eu.

- Lorem ipsum dolor sit amet, consectetur adipiscing elit.
- Vestibulum id tortor efficitur odio dictum suscipit ut blandit augue.
- Donec consectetur sem ac suscipit bibendum.
- Phasellus interdum metus sit amet lorem tempor, vel accumsan mi maximus.
- Suspendisse semper massa et porttitor gravida.

Objetivos y Metodología

Objetivos concretos y metodología de trabajo

[Este bloque es el puente entre el estudio del dominio y la contribución a realizar. Según el tipo concreto de trabajo, el bloque se puede organizar de distintas formas, pero los siguientes elementos deberían estar presentes con mayor o menor detalle.]

2.1. Objetivo General

El objetivo general del presente trabajo es extender la capacidad de monitorización que se tiene sobre una red empresarial. Se enfocará para ello en la detección de anomalías, haciendo uso de técnicas basadas en los equipos que la componen (a diferencia de las basadas en flujos de tráfico). De este modo, se desarrollará un sistema que categorice mediante clustering los equipos finales de la red y que pueda revelar cuándo la actividad de un equipo se desvía de su comportamiento habitual.

Con ello se espera seguir mejorando las prestaciones ofrecidas, especialmente en seguridad, pero también ayudará en la identificación de problemas de configuración, cambios de tendencia en la red y en definitiva cualquier evento inesperado que pudiera repercutir negativamente en la operativa normal de la empresa monitorizada.

2.2. Objetivos Específicos

La consecución de dicho objetivo general vendrá pautada por las siguientes metas más específicas, que se abordarán en este orden:

 Determinar qué características son las más relevantes a la hora de clusterizar la red y extraerlas en tiempo real.

- Establecer unas categorías básicas mediante clasificación no supervisada, que sean verificables con la documentación que se tiene de la red.
- Comprobar que el sistema detecta anomalías si se le proporcionan datos modificados intencionadamente.
- Refinar el algoritmo de clustering empleado, buscando categorías más específicas que no se estuvieran teniendo en cuenta antes.
- Diseñar una forma de visualización adecuada para los resultados.
- Alcanzar un modo de funcionamiento en tiempo real, en coordinación con los demás mecanismos de monitorización presentes y asegurando una precisión razonable (detecciones correctas frente a falsos positivos).

2.3. Metodología de Trabajo

[Sección previsiblemente modificable y ampliable según avance en el desarrollo del trabajo] Con los anteriores objetivos específicos marcados, se hace necesario establecer un marco de trabajo para llevarlos a cabo. En esta sección se estructurará una metodología de trabajo en base a cuatro fases generales, que tienen como fin último alcanzar el objetivo general. Sus enunciados vendrán acompañados de una descripción, explicando con un poco más de detalle cómo se planea ejecutar cada etapa.

Selección de las características más relevantes para la clasificación

A partir de la heterogénea colección de logs de varios firewalls y eventos de antivirus de equipos Windows de la que se dispone, se identificarán los datos que puedan ser interesantes para el clustering. Se tendrán en cuenta distintos aspectos que respondan a las actividades que realiza un equipo informático en una red empresarial, como por ejemplo el uso de determinados servicios, las visitas a sitios web de una serie de categorías, horarios de actividad, interacción con ciertos servidores, etc. A continuación, se analizarán estas características a través de técnicas estadísticas, gracias a las cuales se adquirirá un entendimiento preliminar de la importancia que cada una supone para la posterior clasificación de instancias. Es posible que se logre la reducción de su dimensionalidad, la eliminación de alguna característica redundante o en definitiva alguna simplificación que haga el clustering más sencillo de ejecutar.

Obtención de categorías mediante clustering

Se pasará entonces a aplicar por primera vez algunos algoritmos de clustering sobre estos datos. Inicialmente se elegirán dos algoritmos simples y ampliamente conocidos como son k-means y DB-SCAN. El primero es un algoritmo exclusivo basado en centroides, mientras que el segundo particiona según la densidad de los clusters formados. Se experimentará con diferentes valores para sus parámetros, de forma que puedan compararse resultados. De cara a determinar la efectividad de los algoritmos, se considerarán medidas objetivas que representen su rendimiento. Los métodos típicos para ello valoran características intrínsecas o derivadas como su complejidad, estabilidad y tiempo de computación, así como métricas de validación internas y externas como la silueta de sus clusters o el índice de Rand. Estas medidas son útiles para evaluar la calidad de resultados del algoritmo, pero debe apreciarse que proporcionan un escaso conocimiento sobre lo que contienen los clusters. Además, se emplearán los indicadores comunes con los que se suele definir la bondad de una técnica de clasificación en aprendizaje automático.

Detección de intrusiones con análisis del clustering

Una vez se cuente con una clasificación satisfactoria, se procederá a la siguiente fase, en la cual se pretende detectar anomalías en los datos. En caso de contar con instancias que se hayan identificado como intrusiones mediante otros métodos, se usarán para probar la capacidad de detección del sistema. Si, llegado el momento de hacer este testeo, no se dispone de este tipo de ejemplos, se elaborarán datos que representen distintas clases de anomalías con diversos grados de evidencia. Dado que también se desea detectar otras clases de anomalías más genéricas (no solo relacionadas con seguridad), se incluirán casos como cambios de configuración o de equipos que hayan cambiado de rol en la red.

■ Evaluación en escenario real

Finalmente, se desplegará el sistema desarrollado en un entorno de producción. Se integrará con el framework de alerting empleado y se ofrecerá una representación visual adecuada que permita al analista aprovechar el valor que el sistema aportará a la hora de revisar casos alertados. Además, se automatizará la extracción de características y se prepararán el resto de componentes del sistema para que funcione en tiempo real.

[Se informará del cumplimiento adecuado del Reglamento General de Protección de Datos]

Estado del Arte

Contexto y estado del arte

[Párrafo introductorio del capítulo. Lorem ipsum dolor sit amet, consectetur adipisicing elit. Natus impedit sint cumque, omnis assumenda, molestias corporis repellat, reprehenderit, ullam labore aliquam. Velit ut, ab amet a recusandae, eaque similique alias!]

3.1. Aprendizaje automático en la clasificación de tráfico

El aprendizaje automático puede resumirse como una colección de técnicas de gran potencial para la minería de datos y el descubrimiento de conocimiento (Nguyen y col., 2008). En concreto, a la hora de extraer este conocimiento, estas técnicas son especialmente eficaces buscando y describiendo patrones estructurales útiles en los datos. Además, la ventaja intrínseca de esta disciplina frente a la exploración que pueda efectuar un especialista es que, al poder definirse algorítmicamente el procedimiento para su aplicación, se hace posible implementar sistemas automatizados sobre equipos informáticos.

Un sistema de *machine learning* aprende automáticamente de la experiencia y perfecciona su base de conocimiento, entendiendo "aprender" como la acción de mejorar su rendimiento en el desempeño de una tarea determinada (Herbert, 1983). Esta mejora debe ser cuantificable objetivamente en base a lo que se denomina como medida de rendimiento. El sistema (como cualquier otro agente computacional en la rama de la inteligencia artificial) recibirá estímulos externos (en este caso, conjuntos de datos) y, en base a ellos y al conocimiento que recogen el algoritmo y los resultados de sus ejecuciones previas, producirá una salida que maximizará la medida de rendimiento establecida.

En terminología de *machine learning*, el conjunto de datos que se toma como entrada se compone de instancias. Una *instancia* simboliza a un individuo específico de la población sobre

la que se trabaja. Cada instancia se representa por sus valores en una serie de *características* o atributos, que no son más que medidas sobre aspectos de interés para el escenario en cuestión. Un conjunto de instancias con ciertas características comunes pertenece a una *clase* o concepto. De este modo, se aprende un concepto cuando, dada una instancia, se logra identificar correctamente con qué clase se corresponde. Este aprendizaje implica también que se es capaz tanto de generalizar la aplicación del nombre de la clase a todos los miembros de la misma como de discriminar a los miembros que pertenecen a otra clase.

Atendiendo a la naturaleza de las clases resultantes, los tipos de aprendizaje se dividen en aprendizaje supervisado y no supervisado. El primer tipo es capaz de clasificar nuevas instancias en clases predefinidas. Por el contrario, el segundo clasifica las instancias en clases no definidas con anterioridad. La técnica de aprendizaje no supervisado principal es el clustering o agrupamiento, que se explicará con detalle más adelante.

En los últimos tiempos, el aprendizaje automático se ha venido usando cada vez más en la clasificación de tráfico IP (Dainotti y col., 2012). Las técnicas de clasificación de tráfico siguen mejorando en acierto y eficiencia, pero la proliferación constante de aplicaciones en Internet con comportamientos muy variados, sumado a los incentivos que tienen ciertos agentes para enmascarar algunas aplicaciones y así evitar el filtrado o blogueo en firewalls, son algunas de las razones por las que la clasificación de tráfico permanece como uno de los muchos problemas abiertos de Internet. Han quedado obsoletos métodos clásicos como la identificación de aplicaciones en base a sus puertos conocidos (aquellos registrados por la IANA), que resulta muy simple y rápida pero también poco fiable. En el otro extremo, las técnicas de "Deep Packet Inspection", que analizan en profundidad el funcionamiento de las aplicaciones desde la perspectiva de su uso de los protocolos o buscan datos específicos en paquetes IP para inferir a qué aplicación pertenecen, suponen una alta carga computacional y habitualmente requieren hardware específico. Además, el buen funcionamiento de un clasificador DPI está supeditado a dos condiciones: que pueda inspeccionar el contenido de los paquetes IP y que sepa cómo interpretar la sintaxis de cada aplicación. La primera condición queda comprometida por la estandarización de las conexiones cifradas, mientras que la viabilidad de la segunda se vería restringida por la complejidad de contar con un repositorio completo y constantemente actualizado del formato de los paquetes que puede generar cada aplicación. Ante estas técnicas también se presentan dificultades legales y relacionadas con la privacidad.

Si se trata de clasificadores de tráfico, lo más común es encontrar planteamientos basados en flujos de tráfico. En ocasiones, la granularidad de la clasificación se afina hasta el uso de

flujos bidireccionales (asumiendo que se tiene visibilidad de ambas direcciones), pero operar a este nivel entraña una complejidad bastante mayor. Un flujo se suele definir como una tupla de 5 elementos: protocolo de transporte (frecuentemente, TCP o UDP), direcciones IP de origen y destino, y puertos de origen y destino. Con este concepto como *objeto* fundamental, tradicionalmente se han buscado patrones estadísticos en los atributos de los flujos que son observables desde una perspectiva externa (es decir, sin considerar el contenido o *payload* de los paquetes). Ejemplos de estos atributos serían: tamaño de paquetes, tiempo entre llegadas, número de paquetes en cada dirección, duración total... resumido cada uno con el estadístico muestral que se considere adecuado. Con la popularización del aprendizaje automático, se ha podido llevar la búsqueda de patrones entre dichos atributos a nuevos grados de profundidad.

Se trabaja también con otras variantes en cuanto a cómo agrupar los paquetes que se hayan intercambiado dos máquinas. Entre ellas, podrían destacarse las conexiones TCP o los servicios, definidos estos como el tráfico generado entre una pareja de IPs-puertos. En cualquiera de los casos anteriores, se pone el foco sobre flujos individuales, para después clasificarlos bajo categorías que comparten características. Este tipo de planteamientos no tienen tan en cuenta el conjunto de acciones que lleva a cabo un mismo equipo. Así, se corre el riesgo de perder información útil de cara a entender de forma completa qué es realmente lo que están haciendo los equipos de la red.

Por otro lado, se encuentran los clasificadores de tráfico basados en el comportamiento del host. Sirva de referente el trabajo de (Karagiannis y col., 2005), donde se propuso un novedoso método que identificaba patrones en el comportamiento de los hosts a la altura de la capa de transporte. Se trataba de una aproximación multinivel que descartaba incluir datos sobre el payload, los puertos bien conocidos (aspectos que conllevan las problemáticas anteriormente comentadas) o cualquier otra información separada de la que ofrecían los colectores de flujos. Consistía por tanto en un clasificador "a ciegas" ("BLINd Classification", abreviado como "BLINC") que analizaba cada hosts desde tres perspectivas: social, funcional y aplicativa. La perspectiva social capturaba las interacciones del host con otros hosts, en términos de cuántos hosts se conectaban con qué hosts. La funcional los separaba según actuaran como proveedores de un servicio, consumidores o ambos. Se tenían en cuenta, por tanto, los roles del modelo cliente-servidor. Por último, en la perspectiva aplicativa se utilizaba la información de la capa de transporte con la intención de distinguir la aplicación en cuestión.

Mediante la premisa de no tratar cada flujo como una entidad distinta, se conseguiría acumular la información necesaria para reconocer el verdadero comportamiento de cada host. Además de cumplir con la identificación de aplicaciones específicas, este método sería

resistente a circunstancias de la red como congestión o cambios de rutas. Esto es así porque, a diferencia de otros métodos (véanse los mencionados sobre flujos), una aproximación centrada en el comportamiento de los hosts suele ser insensible a las variaciones que puedan presentar parámetros como los tiempos de llegada entre paquetes. En cuanto a resultados, los enfoques de aprendizaje automático basados en patrones de comunicación de los hosts alcanzan resultados comparables a los de técnicas de DPI, siendo notablemente más asequibles y menos invasivos con la privacidad.

Es por todo lo anterior que en los sistemas de detección de intrusiones, que se van a desarrollar en la siguiente sección, priman los enfoques sobre el host en vez de sobre el flujo.

3.2. Detección de anomalías sobre actividad de red

En la gestión y monitorización de una red empresarial cobra especial relevancia la seguridad. En este ámbito, la seguridad informática se centra en proteger la red corporativa de ataques que puedan comprometer su disponibilidad o la integridad de los equipos que la componen, así como bloquear acciones no autorizadas y evitar el uso indebido de los recursos que quedan expuestos al exterior.

Las organizaciones toman numerosas medidas de seguridad frente a estas amenazas, tanto *software* como *hardware*. Dos ejemplos claros serían los antivirus y los firewalls, que podríamos englobar dentro de las aproximaciones a la seguridad basadas en firmas (D'Alconzo y col., 2019). Sin embargo, estos métodos dependen de que el fabricante del producto de seguridad haya detectado el ataque previamente, haya generado una firma que lo identifique y haya distribuido la misma hasta el cliente final. Es decir, solo pueden ofrecer protección ante ataques conocidos y requieren que todos los pasos anteriores se hayan completado.

En contraposición, existen los sistemas de seguridad basados en detección de anomalías. Este tipo de métodos asumen que el impacto de un ataque modificará el comportamiento de la red, así que construyen un modelo que represente el comportamiento normal de la red, especificado por ciertas métricas. A continuación, monitorizan el tráfico y fijan alarmas que se dispararán cuando el valor recogido en alguna de esas métricas de referencia se desvíe del rango considerado normal (Boutaba y col., 2018).

Habitualmente, este tipo de defensas basadas en detección de anomalías son complementarias a las basadas en firmas. Se sitúan en una segunda línea con el objetivo de detectar a tiempo síntomas tempranos de ciberataques, para así poder actuar antes de que causen daños. Ambos enfoques pueden encontrarse integrados en soluciones conocidas como IDS/IPS (Intrusion Detecion/Prevention System).

Hablando en términos generales, pueden distinguirse tres fases básicas que cumplen todos los NIDS (*Network* IDS) basados en anomalías (García-Teodoro y col., 2009):

- Parametrización: las instancias del sistema objetivo se representan de forma adecuada para su tratamiento.
- Entrenamiento: se caracteriza el comportamiento normal del sistema, mediante un modelo que puede construirse con técnicas basadas en alguna de las categorías descritas después.
- Detección: se compara el modelo con el tráfico disponible, de forma que se dispara una alarma si una instancia se desvía.

Según cómo se modele el comportamiento normal del sistema (Lazarevic y col., 2005), se pueden categorizar en técnicas basadas en estadística, en conocimiento o en aprendizaje automático. Las primeras no requieren un conocimiento previo sobre la actividad normal del sistema, pero la presunción que asumen de cuasi-estacionalidad es poco realista. Las segundas son robustas, pero el mantenimiento de datos de calidad resulta difícil y costoso. En cuanto a las técnicas basadas en aprendizaje automático, son flexibles y adaptables. También pueden capturar interdependencias entre las variables que no son fáciles de encontrar de otra forma. No obstante, estas técnicas tienen una dependencia importante de lo que se acepte como comportamiento normal.

En (D'Alconzo y col., 2019) se resaltan acertadamente las bondades y debilidades de los métodos de detección de anomalías, al decir que "son atractivos porque permiten la pronta detección de amenazas desconocidas (por ejemplo, zero-days). Estos métodos, sin embargo, puede que no detecten ataques sigilosos, insuficientemente amplios para perturbar la red. A veces también adolecen de un alto número de falsos positivos."

Continúa señalando cómo beneficia el aprendizaje automático a este tipo de sistemas: "el machine learning ha recibido una significativa atención en la detección de anomalías, debido a la autonomía y robustez que ofrece en el aprendizaje y también a la hora de adaptar el perfil de la normalidad según va cambiando. Con machine learning, el sistema puede aprender patrones de comportamientos normales dentro de entornos, aplicaciones, grupos de usuarios y a lo largo del tiempo. Además, ofrece la capacidad de encontrar correlaciones complejas en los datos que no pueden deducirse de la mera observación".

Se concluye por tanto que la obtención de una representación completa de la normalidad, requisito no trivial en estos sistemas basados en detección de anomalías, puede tomarse como un problema de clasificación en instancias normales y no normales. Dicho problema puede

abordarse mediante la técnica de aprendizaje no supervisado descrita a continuación.

3.3. Clustering

El clustering se define en (Nguyen y col., 2008) como la agrupación de instancias que tienen características *cercanas* en forma de clusters, sin aplicar ninguna orientación previa. Esta técnica de aprendizaje automático no supervisado asocia a las instancias con propiedades similares bajo el mismo grupo, determinando dicha similaridad en un modelo que posibilite la medición de distancias específicas, como pueda ser el espacio euclídeo. Los grupos pueden ser exclusivos, si cada instancia pertenece a un único grupo; solapados, si una instancia puede pertenecer a varios grupos; o probabilísticos, si la pertenencia de una instancia a un grupo se expresa mediante una cierta probabilidad.

El primer uso de clustering para detección de intrusiones se vio en (Portnoy, 2000). La hipótesis en base a la cual los autores aplicaron clustering para esta tarea es que las conexiones entre datos normales crearán clusters más grandes y más densos. Si se lleva el análisis un paso más allá, para incrementar la precisión de la técnica, además de las consideraciones anteriores debe tenerse en cuenta la distancia entre clusters, como se demuestra en (Jiang y col., 2006).

Los tipos de algoritmos de clustering usados en clasificación de tráfico son variados. En (McGregor y col., 2004) se aplica por primera vez un algoritmo de clustering probabilístico como es *Expectation Maximization* sobre flujos de tráfico. Considerando varias estadísticas sobre longitud de paquetes, tiempos entre llegadas, cantidad de bytes, duración de la conexión y número de permutaciones de modo "transaccional" a modo "por lotes" (y viceversa), se consigue una clasificación con baja granularidad. Otros estudios valoran más medidas estadísticas [Más algoritmos de clustering en clasif de tráfico: - AutoClass (Zander et al.): - features: packet length stats (mean, variance in forward and backward directions), inter-arrival stats (mean, var. in f-b dirs), flow size (bytes), flow duration. Calculated on full flows - classif level: (muy granular) 8 aplicaciones estudiadas - K-Means (Bernaille et al.): - features: Packet lengths of the first few packets of bi-directional traffic flows - features: (muy granular) 10 aplicaciones estudiadas]

```
[Algoritmos de clustering en detección de anomalías: (Syarif y col., 2012) (Leung y col., 2005) (Kim y col., 2018) ]
```

```
[Network-IDS basados en anomalías: (Bhuyan y col., 2014)]
```

[Extender las ideas y las referencias de: (Bohara y col., 2016)]

[Conclusiones sobre los trabajos previos]

Desarrollo

Desarrollo específico de la contribución

[Párrafo introductorio del capítulo. Lorem ipsum dolor sit amet, consectetur adipisicing elit. Natus impedit sint cumque, omnis assumenda, molestias corporis repellat, reprehenderit, ullam labore aliquam. Velit ut, ab amet a recusandae, eaque similique alias!]

4.1. Presentación del entorno

En el escenario del proyecto (la red de un banco que es cliente de la empresa), [..]

4.2. Extracción y filtrado

El punto de interés para la captura se concentra por tanto en dos equipos por los que pasa el grueso del tráfico total: un Firewall Fortinet y un IPS PaloAlto. Como es habitual, estos equipos reportan todas sus acciones a través de logs, con varios niveles de granularidad e importancia. Incluyen también multitud de información aportada por los propios sistemas que enriquecen el valor de cada evento. Esto es razón para preferir los logs de firewall como fuente de información frente a una captura de tráfico en crudo, ya que lo que procesan los firewalls casi siempre es más relevante que el tráfico completo pero sin procesar. Además, el volumen de una captura de tráfico de estas características sería notablemente mayor y más difícil de manejar.

Aunque el fin para el que sirven ambos equipos (análisis y protección frente a amenazas informáticas) sea similar, la estructura usada por cada uno en los logs que produce es completamente distinta. Los logs de Fortinet siguen un formato clave-valor con ciertas particularidades, mientras que los de PaloAlto tienen una serie de campos fijos que están delimitados por comas. A modo de ejemplo, las dos líneas adjuntadas a continuación corresponden a un evento del firewall Fortinet y otro del IPS PaloAlto, respectivamente (cada evento viene en una única línea):

```
1585572524|1585572524|2020-03-30T06:48:44.202297|10.2.0.11|6|local7|

date=2020-03-30 time=06:48:44 devname="FW1_INTERNETCORP" devid="FG1809999"

logid="1059028704" type="utm" subtype="app-ctrl" eventtype="app-ctrl-all"

level="information" vd="root" eventtime=1585572524 appid=41470 user="NOM"

group="GrupoOffice365" authserver="SV1" srcip=172.2.9.6 dstip=23.203.51.72

srcport=54697 dstport=443 srcintf="p18" srcintfrole="undef" dstintf="p20"

dstintfrole="wan" proto=6 service="HTTPS" direction="outgoing" policyid=124

sessionid=325186437 applist="AC_CORREO" appcat="Collab" app="Microsoft.CDN"

action="pass" hostname="img-prod-cms-rt-microsoft-com.akamaized.net"

incidentserialno=1513678724 url="/" msg="Collaboration: Microsoft.CDN,"

apprisk="elevated" scertcname="a248.e.akamai.net"
```

1585659863|1585659863|2020-03-31T07:04:23.027791|10.2.0.73|6|local0|
1,2020/03/31 07:04:23,001801037558,TRAFFIC,end,2049,2020/03/31 07:04:03,
10.138.4.7,186.151.236.155,0.0.0.0,0.0.0,0UTBOUND,,,incomplete,vsys1,
trust,untrust,ethernet1/10,ethernet1/9,Log-Panorama,2020/03/31 07:04:03,
41602,1,55074,80,0,0,0x19,tcp,allow,132,132,0,2,2020/03/31 07:03:55,3,any,
0,1307298109,0x80000,10.0.0.0-10.255.255.255,America,0,2,0,aged-out,13,0,0,0,0,PA-3020-Z9,from-policy,,,0,,0,,N/A,0,0,0,0

[Volumen de estos logs]

[Gráfica de volumen acumulado en el tiempo]

Otro hecho reseñable que afecta al formato es que se emplea syslog¹ (el estándar de facto) como protocolo para trasladar los datos desde cada equipo hasta la sonda, de forma que se cuenta con ciertos campos adicionales a los enviados por los equipos. Para el tema que nos ocupa, los únicos campos que se extraen de esta cabecera son: la marca de tiempo en la que ha llegado cada evento, conocida en el vocabulario informático como timestamp, y la prioridad del evento. En cualquier caso, esta sección adicional dentro de los logs tiene también su propio formato, por lo cual también se deberá tratar de forma específica. En nuestra

https://tools.ietf.org/html/rfc5424

configuración (que aplica a la herramienta *rsyslog*²), la siguiente directiva establece cómo se vuelcan a fichero estos campos de syslog:

Así que, en los *scripts* que procesan los ficheros donde se han volcado los datos traídos mediante *syslog*, se obtiene la fecha de cada evento a partir de este primer campo "timereported:::date-unixtimestamp" y la prioridad a partir del quinto campo. Esta primera parte del procesado (que está programado en Python) se hace de la siguiente manera:

```
for syslogline in sys.stdin:
    try:
        splitted_syslogline = syslogline.rstrip().split("|")
        → #.rstrip() removes last "\n" character
        tstamp_line = int(splitted_syslogline[0])
        prio = splitted_syslogline[4]
```

Seguidamente, el resto de la línea actual (sin la cabecera de *syslog*) se convierte en una estructura de diccionario. Como se apreciaba en las líneas de ejemplo que se han incluido antes, la relación entre claves y valores depende de cada caso.

Para el equipo Fortinet, la relación está definida en el propio evento como clave="valor" o clave=valor para valores no considerados como cadenas de texto. Cabe destacar que el símbolo "=" puede estar contenido en el valor. El nombre de la clave, sin embargo, nunca llevará comillas. Establecidas las anteriores reglas, en Fortinet se convierten los campos con una expresión regular y una dict comprehension³ (es decir, una forma concisa de crear diccionarios a través de la iteración sobre una lista con la posibilidad de incluir condicionales):

https://www.rsyslog.com/

³https://www.python.org/dev/peps/pep-0274/

Para el IPS PaloAlto, la extracción de los campos es más sencilla. Como cumplen con el formato CSV estandarizado⁴, los valores se tienen en una lista con solo leer la línea a través de una función de la librería CSV. En cuanto a las claves, en la documentación⁵ de PaloAlto se explica que depende del tipo del evento. Por tanto, se asignan unas claves u otras consultando primero de qué tipo se trata. Finalmente, se construye el diccionario con otra *dict comprehension*:

Posteriormente se lleva a cabo otra serie de operaciones necesarias para la transformación de los datos de entrada en información útil para la monitorización. Sin embargo, desde la perspectiva de este trabajo, el único apartado de interés es la agregación de sesiones, que se desarrolla a continuación.

Una parte de este procesado consiste en guardar cierta información asociada a cada sesión. El concepto de "sesión" sería equivalente al de "flujo" presentado en el capítulo anterior: una serie de eventos asociados que se corresponden con la misma tupla de {IP origen, IP destino, protocolo, puerto origen, puerto destino}. Los dos equipos mantienen un campo "Identificador de Sesión" interno que se añade a la tupla de la sesión. Este campo permite distinguir los eventos de sesiones que coinciden en origen y destino pero se producen en intervalos temporales diferentes. Se incluye en el procesado con esta finalidad de no confundir sesiones distintas en etapas posteriores, pero para nuestro propósito de clasificación de equipos puede ignorarse.

La información de cada sesión se compone de:

■ Tupla que define la sesión:

{ID de sesión, IP origen, IP destino, protocolo, puerto origen, puerto destino}

⁴https://tools.ietf.org/html/rfc4180

⁵https://docs.paloaltonetworks.com/pan-os/8-1/pan-os-admin/monitoring/ use-syslog-for-monitoring/syslog-field-descriptions.html

- Timestamps del primer y último evento pertenecientes a esta sesión
- Máxima prioridad de evento vista en esta sesión
- Bytes recibidos y enviados (solo en el IPS)
- Nivel de anomalía
- Nivel de amenaza
- Contador y lista de eventos

Los niveles de anomalía y amenaza son unos índices simples que se han diseñado para resumir cualidades de interés acerca de la sesión, como son cuánto se aleja de la normalidad la cantidad de eventos prioritarios que se han visto y cómo son de peligrosas las amenazas recibidas. Para cada sesión, se calculan sus niveles con las fórmulas:

$$N_{\text{anomalia}} = \sum_{i=1}^{n} \frac{1}{\text{prioridad}_{\text{evento}_i}}$$

para eventos de prioridad ≤ 4 o eventos de tráfico que no son de inicio ni fin

$$N_{\mathrm{amenaza}} = \sum_{i=1}^{n} \frac{1}{\mathrm{prioridad}_{\mathrm{evento}_i}}$$

para eventos de amenaza con prioridad ≤ 4 que no son bloqueados

Tanto estas como el resto de características cuantificables pueden resultar de interés a la hora de aplicar *clustering* sobre un dataset derivado de este procesado.

La recolección de esta información se realiza a través de la función adjunta. Como puede verse, cuando se llama a esta función (lo cual ocurre ante todos los eventos de protocolo TCP o UDP que incluyan ID de sesión) se actualizan los parámetros relativos a la sesión actual, que están almacenados en un diccionario. A su vez, este diccionario se encuentra dentro del diccionario sessions. Con él se mantienen en memoria todas las sesiones que todavía no se han cerrado. Cuando transcurre un *bucket* de tiempo determinado (por defecto, 60 segundos), se comprueban todas las sesiones vigentes. Aquellas que tienen un evento de finalización se imprimen en un fichero y se retiran del diccionario sessions. De este modo, cada 60 segundos se tienen nuevas sesiones completas (en la práctica se alcanzan incluso más de 100 000 sesiones finalizadas cada minuto) que serán datos de entrada para el *clustering*.

```
def session_aggregation(dict_line, event_descript, event_tstamp):
  It groups info related to actual session on the sessions dictionary, that is:
  - the sessionid and the session tuple (srcip-dstip-proto-srcport-dstport)
 - tstamp of first and last event observed for actual session

    update max. priority of events observed for actual session

  - sent and received bytes for actual session
  - counter and list of events observed for actual session
  - recalculate anomaly level for actual session
  - recalculate threat level for actual session
    session_tuple = "...".join([
             dict_line['SESSION ID'], dict_line['SRC_IP'], dict_line['DST_IP'],
             dict_line['PROTO'], dict_line['SRC_PORT'], dict_line['DST_PORT']
        ])
    priority = int(dict_line['priority'])
    if session_tuple not in sessions:
        sessions[session_tuple] = {}
        # store the session tuple values related to this new session_tuple:
        sessions[session_tuple]['events'] = []
        sessions[session_tuple]['anomaly_level'] = 0
        sessions[session_tuple]['threat_level'] = 0
        sessions[session_tuple]["counter"] = 0
        sessions[session_tuple]['max_prio'] = priority
        sessions[session_tuple]["bytes_sent"] = 0
        sessions[session_tuple]["bytes_rcvd"] = 0
        sessions[session_tuple]['first_event_tstamp'] = int(event_tstamp)
    sessions[session_tuple]['last_event_tstamp'] = int(event_tstamp)
    sessions[session_tuple]["counter"] += 1
    if dict_line['type'] == "TRAFFIC":
         sessions[session_tuple]["bytes_sent"] += int(dict_line['BYTES_SENT'])
        sessions[session_tuple]["bytes_rcvd"] += int(dict_line['BYTES_RECEIVED'])
    if priority < sessions[session_tuple]['max_prio'];</pre>
    # lower prio value means more important (i.e., the most important priority is 1, or even 0 if priority=0 exists)
        sessions[session_tuple]['max_prio'] = priority
    else:
        sessions[session_tuple]['max_prio']
    if priority<=4 or (dict_line['type']=="TRAFFIC" and "end" not in event_descript</pre>
     → and "start" not in event_descript):
        sessions[session_tuple]['anomaly_level'] += 1/priority
    if priority<=4 and dict_line['type']=="THREAT" and dict_line['ACTION']!="alert":</pre>
        sessions[session_tuple]['threat_level'] += 1/priority
    sessions[session_tuple]['events'].append("{}, {}".format(event_descript,

    dict_line['ACTION']))
```

4.3. Preprocesado para el clustering

4.4. Análisis de datos

Una vez extraídas las características, se procede a analizar la distribución de valores que presenta cada una de ellas en este *dataset*. Dicha tarea permitirá entender mejor la importancia relativa de cada característica y cómo afectará a los algoritmos de *clustering*. En concreto, se prestará especial atención a la forma, varianza y modalidades de las distribuciones.

[Funciones de distribución acumulada para cada característica]

4.5. Selección de características

Como se apunta en "Análisis de las características en tráfico de red para detección de anomalías" (Iglesias y col., 2015), la meta que tiene la selección de características en la detección de anomalías es "eliminar características fuertemente correladas, relevantes e irrelevantes para mejorar la calidad de la detección". En este trabajo, los autores abordaron la selección de características de forma exhaustiva y rigurosa, mediante métodos multi-fase implementados con envolvedores (*wrappers*, que buscan el subconjunto de características con mejores resultados), combinando filtrado y técnicas de regresión gradual.

[Diferencia entre selección de feats.]

También se considera la idea de (Guyon y col., 2003) cuando dice que "el objetivo de la selección de variables es triple: mejorar el rendimiento de los predictores, posibilitar predictores más rápidos y eficientes en coste, y permitir una mejor comprensión del proceso subyacente que ha generado los datos".

4.6. Parametrización

Resultados

Explicación de los resultados

[Párrafo introductorio del capítulo. Lorem ipsum dolor sit amet, consectetur adipisicing elit. Suscipit maxime expedita possimus consequatur labore, id, dignissimos praesentium repellendus quisquam tempore natus eveniet magnam! Ad laborum quas, at tenetur eligendi officiis!]

Conclusiones

Conclusiones finales y trabajo futuro

[Párrafo introductorio del capítulo. Lorem ipsum dolor sit amet, consectetur adipisicing elit. Suscipit maxime expedita possimus consequatur labore, id, dignissimos praesentium repellendus quisquam tempore natus eveniet magnam! Ad laborum quas, at tenetur eligendi officiis!]

Este último bloque (habitualmente un capítulo; en ocasiones, dos capítulos complementarios) es habitual en todos los tipos de trabajos y presenta el resumen final de tu trabajo y debe servir para informar del alcance y relevancia de tu aportación.

Suele estructurarse empezando con un resumen del problema tratado, de cómo se ha abordado y de por qué la solución sería válida. Es recomendable que incluya también un resumen de las contribuciones del trabajo, en el que relaciones las contribuciones y los resultados obtenidos con los objetivos que habías planteado para el trabajo, discutiendo hasta qué punto has conseguido resolver los objetivos planteados.

Finalmente, se suele dedicar una última sección a hablar de líneas de trabajo futuro que podrían aportar valor añadido al TFM realizado. La sección debería señalar las perspectivas de futuro que abre el trabajo desarrollado para el campo de estudio definido. En el fondo, debes justificar de qué modo puede emplearse la aportación que has desarrollado y en qué campos.

Bibliografía

- Herbert, S (1983). «Why Should Machines Learn?» En: *Machine Learning: An Artificial Intelligence Approach*. Springer Berlin Heidelberg, págs. 25-37. isbn: 978-3-662-12405-5. doi: 10.1007/978-3-662-12405-5_2.
- Portnoy, L. (2000). «Intrusion Detection with Unlabeled Data Using Clustering». En: doi: 10.7916/D8MP5904.
- Guyon, I. y A. Elisseeff (2003). «An Introduction to Variable and Feature Selection». En: *J. Mach. Learn. Res.* 3, págs. 1157-1182. doi: 10.5555/944919.944968.
- McGregor, A. y col. (2004). «Flow Clustering Using Machine Learning Techniques». En: *Passive and Active Network Measurement* 3015, págs. 205-214. doi: 10.1007/978-3-540-24668-8 21.
- Karagiannis, T., K. Papagiannaki y M. Faloutsos (2005). «BLINC: Multilevel Traffic Classification in the Dark». En: *SIGCOMM Computer Communications Review 35.4*, págs. 229-240. doi: 10.1145/1090191.1080119.
- Lazarevic, A., V. Kumar y J. Srivastava (2005). «Intrusion Detection: A Survey». En: vol. 5, págs. 19-78. doi: 10.1007/0-387-24230-9_2.
- Leung, K. y C. Leckie (2005). «Unsupervised Anomaly Detection in Network Intrusion Detection Using Clusters». En: págs. 333-342. isbn: 1920682201.
- Jiang, S. y col. (2006). «A clustering-based method for unsupervised intrusion detections». En: *Pattern Recognition Letters* 27, págs. 802-810. doi: 10.1016/j.patrec.2005.11.007.
- Nguyen, H.T. y G. Armitage (2008). «A survey of techniques for internet traffic classification using machine learning». En: *IEEE Communications Surveys & Tutorials* 10.4, págs. 56-76. doi: 10.1109/SURV.2008.080406.
- García-Teodoro, P. y col. (2009). «Anomaly-based network intrusion detection: Techniques, systems and challenges». En: *Computers & Security* 28, págs. 18-28. doi: 10.1016/j.cose.2008.08.003.
- Dainotti, A., A. Pescape y K. C. Claffy (2012). «Issues and future directions in traffic classification». En: *IEEE Network* 26.1, págs. 35-40. doi: 10.1109/MNET.2012.6135854.

- Syarif, I., A. Prugel-Bennett y G. Wills (2012). «Unsupervised clustering approach for network anomaly detection». En: doi: 10.1007/978-3-642-30507-8_7.
- Bhuyan, M. H., D. K. Bhattacharyya y J. K. Kalita (2014). «Network Anomaly Detection: Methods, Systems and Tools». En: *IEEE Communications Surveys & Tutorials* 16.1, págs. 303-336. doi: 10.1109/SURV.2013.052213.00046.
- Iglesias, F. y T. Zseby (2015). «Analysis of network traffic features for anomaly detection». En: *Machine Learning*, págs. 59-84. doi: 10.1007/s10994-014-5473-9.
- Bohara, A., U. Thakore y W. Sanders (2016). «Intrusion Detection in Enterprise Systems by Combining and Clustering Diverse Monitor Data». En: *Proceedings of the Symposium and Bootcamp on the Science of Security*, págs. 7-16. isbn: 9781450342773. doi: 10.1145/2898375.2898400.
- Boutaba, R. y col. (2018). «A Comprehensive Survey on Machine Learning for Networking: Evolution, Applications and Research Opportunities». En: *Journal of Internet Services and Applications* 9. doi: 10.1186/s13174-018-0087-2.
- Kim, J. y col. (2018). «Multivariate network traffic analysis using clustered patterns». En: *Computing* 101. doi: 10.1007/s00607-018-0619-4.
- D'Alconzo, A. y col. (2019). «A Survey on Big Data for Network Traffic Monitoring and Analysis». En: *IEEE Transactions on Network and Service Management* 16.3, págs. 800-813. doi: 10.1109/TNSM.2019.2933358.

Apéndices

Apéndice A

Fusce viverra lectus

Nam viverra, odio et vulputate ultricies, sem libero ornare nunc, nec suscipit urna sapien eu sapien. Nulla erat nulla, hendrerit non elit vitae, venenatis malesuada libero. Duis ornare sapien sed lacus condimentum rutrum. Donec feugiat erat id elit aliquam, a ultrices lectus mattis. In vitae nisl tortor. Proin pellentesque nec odio et posuere. Ut aliquet quam ac magna tincidunt ultrices. Nullam sit amet elementum leo. Pellentesque vitae mi dolor. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Suspendisse tincidunt mi arcu, vitae porttitor mauris elementum ut.