



UNIVERSITAT OBERTA DE CATALUNYA (UOC)
MÁSTER UNIVERSITARIO EN CIENCIA DE DATOS (*Data Science*)

TRABAJO FINAL DE MÁSTER

ÁREA: MINERÍA DE DATOS Y MACHINE LEARNING

Detección de anomalías en entorno del Internet de las cosas

Autor: Gonzalo Pedro Mellizo-Soto Díaz

Tutor: Carlos Hernández Gañán

Profesor: Jordi Casas Roma

Madrid, 5 de junio de 2019

Copyright



Esta obra está sujeta a una licencia de Reconocimiento - NoComercial - SinObraDerivada

3.0 España de Creative Commons.

FICHA DEL TRABAJO FINAL

Título del trabajo:	Detección de anomalías en el entorno del Internet de las cosas
Nombre del autor:	Gonzalo Pedro Mellizo-Soto Díaz
Nombre del colaborador/a docente:	Carlos Hernández Gañán
Nombre del PRA:	Jordi Casas Roma
Fecha de entrega (mm/aaaa):	06/2019
Titulación o programa:	Máster Universitario en Ciencia de Datos
Área del Trabajo Final:	Minería de datos y Machine Learning
Idioma del trabajo:	Español
Palabras clave	Machine Learning, IOT, Anomaly Detection

Cita

“Nuestro lema es: más humanos que los humanos”

Eldon Tyrell, *Blade Runner*

Abstract

In recent years the amount of connected devices has greatly increased, with an increasing number of applications in the industry each day. These devices can be subject of attacks causing instability or data leaks that can be dangerous both for the users and the enterprises, in order to avoid or confront them, security and early detection are becoming a must in a connected world. The focus is the monitoring and detection of the attacks in Internet of Things devices using state of the art Machine Learning techniques. Models such as SVM, DBScan or Isolation Forests have been used and assembled in order to identify with a better accuracy when an attack is happening. With this assembly, attack detection has increased up to 15 % comparing to traditional methods and individual model usages and times have been considerably reduced. An active use of Machine Learning models has shown a great improvement at anomaly detection by securing the devices and decreasing the reaction times when facing attacks.

Durante los últimos años se encuentra una creciente cantidad de dispositivos conectados entre sí, cada vez con más aplicaciones en la industria. Estos dispositivos pueden ser atacados y provocar inestabilidad o una fuga de datos, por lo tanto la protección y la pronta detección de ataques y/o anomalías es vital en un mundo cada vez más conectado. El objetivo es la monitorización y detección de estos ataques en dispositivos del *Internet of Things* utilizando técnicas del estado del arte de Machine Learning para su detección y poder así responder con una mayor rapidez a los ataques. Para la detección se han utilizado modelos estadísticos, como SVM, DBScan o Isolation Forests, que en su conjunto permitan identificar con mayor precisión cuando se está produciendo un ataque. El conjunto de la clusterización con la clasificación de puntos anómalos muestra una mayor robustez, frente al uso individual de cada uno de los modelos aumentando la detección en hasta un 15 %. Se demuestra cómo el uso de los modelos permite proteger los dispositivos y mejorar la seguridad al disminuir los tiempos de reacción frente a los ataques.

Palabras clave: Machine Learning, IOT, Anomaly Detection

Índice general

Abstract	VII
Índice	IX
Listado de Figuras	XI
Listado de Tablas	1
1. Introducción	3
1.1. Contexto y justificación del Trabajo	3
1.2. Explicación de la motivación personal	4
1.3. Objetivos del Trabajo	4
1.4. Descripción general del problema	5
1.5. Enfoque y método seguido	5
1.6. Planificación del Trabajo	5
2. Estado del Arte	9
2.1. Métodos tradicionales de detección de anomalías	10

2.2. Machine Learning y la detección anomalías	11
2.2.1. Aprendizaje Supervisado	13
2.2.2. Aprendizaje No Supervisado	22
3. Desarrollo	29
3.1. Datos utilizados	30
3.2. Exploración y Preprocesamiento	33
Bibliografía	42

Índice de figuras

1.1. Planificación de tareas	7
2.1. Relación entre Inteligencia Artificial, Machine Learning y Deep Learning	12
2.2. Regresión Lineal	14
2.3. Red Neuronal Artificial	15
2.4. Arquitectura de un Perceptron	15
2.5. Función de activación: Función Escalón	16
2.6. Backpropagation	17
2.7. Hiperplanos SVM	20
2.8. Hiperplanos para distintos valores de C	21
2.9. Efecto de aplicación de <i>Kernel</i> radial en SVM	21
2.10. Datos de dos variables sin PCA	23
2.11. PCA aplicado a dos variables	23
2.12. Ejemplo de agrupamiento	25
2.13. Arquitectura de un Autoencoder	27

2.14. Ejemplo de detección de anomalías con Autoencoders	28
3.1. Ejemplo de dirección IPv4	31
3.2. Ejemplo de datos pcap en Wireshark	32
3.3. Datos en forma tabular	35
3.4. Cantidad de valores únicos en las variables	35
3.5. Puertos más utilizados	36
3.6. IPs de origen con la mayor cantidad de conexiones	37
3.7. IPs de destino con la mayor cantidad de conexiones	37
3.8. Evolución de las conexiones a lo largo del tiempo con distintas ventanas	38
3.9. Histograma de la variable <i>Length</i>	39

Índice de cuadros

3.1. Tabla de frecuencias de las variables binarias	39
---------------------------------------------------------------	----

Capítulo 1

Introducción

1.1. Contexto y justificación del Trabajo

Cada vez se encuentran más dispositivos conectados entre sí, no solo en la industria, sino también en los hogares, esta conexión entre dispositivos los hace vulnerables a ataques informáticos que pueden afectar en gran medida a los usuarios, no solo pueden provocar un mal funcionamiento de los mismos, sino que también puede provocar la fuga de datos de distinta sensibilidad. La previsión en los futuros años es estar cada vez más conectados y una pronta detección de los ataques puede ayudar a evitar los problemas derivados de los mismos, mediante una pronta reacción o haciendo frente a la previsión de un ataque.

Actualmente se utilizan distintas medidas de seguridad como control de acceso físico al dispositivo, encriptación de datos, firewalls, securización de red, etc... Sin embargo, en muchos casos la detección de anomalías aplica un umbral estacionario, provocando que en el caso de un ataque ya sea demasiado tarde para reaccionar o no se sea capaz de identificar patrones extraños previos al ataque. Con la aplicación de nuevas técnicas se pretende mejorar los tiempos de reacción e incluso predecir cuándo puede suceder un ataque.

1.2. Explicación de la motivación personal

La razón principal de la elección del proyecto es la posibilidad de aprender y utilizar técnicas de Machine Learning en un sector desconocido que permita diversificar conocimientos. Todo se encuentra cada vez más conectado e investigar cómo clasificar cuando se está produciendo un ataque permite profundizar en conocimientos de seguridad y aplicarlos en un problema real.

La aplicación de técnicas en un problema real ayuda a comprender mejor el uso de las herramientas y el porqué y cuándo se deben de utilizar. De este modo, se añade un nuevo conocimiento que aportar al ámbito profesional y puede utilizarse también en un entorno privado.

1.3. Objetivos del Trabajo

Los objetivos del trabajo son los siguientes:

- Adquisición de conocimiento del sector y de los dispositivos IOT.
- Lectura y comprensión de las técnicas del estado del arte en detección de anomalías en dispositivos conectados.
- Obtención de datos reales de conexiones a dispositivos IOT, en caso de no ser posible, generación de datos sintéticos.
- Prueba de las técnicas encontradas y evaluación en el problema actual.
- Investigación de algoritmos tradicionales de Machine Learning y su utilidad.
- Desarrollo de la solución utilizando los algoritmos más aptos.
- Evaluación de la detección de patrones y ataques de las técnicas utilizadas.
- Comparación de los resultados obtenidos con el estado del arte.

1.4. Descripción general del problema

La cantidad de dispositivos informáticos existentes que pueden ser víctimas de un ataque es masiva, por lo que securizar correctamente los dispositivos es fundamental. A pesar del uso de distintos protocolos de seguridad, se generan nuevos tipos de ataque todos los días, por lo tanto la detección es una necesidad a la hora de evitar los problemas derivados.

La pérdida de control de los dispositivos o la fuga de información de los mismos puede provocar pérdidas millonarias a las empresas o provocar una gran inseguridad a los usuarios de productos IOT, así como causar posibles daños a infraestructuras o personas. Por otro lado, muchos de los dispositivos pueden no tener la capacidad de computación necesaria para incluir una capa de seguridad robusta, que puede compensarse con una detección temprana.

1.5. Enfoque y método seguido

El enfoque pasa por conocer el problema y dar respuesta a preguntas específicas, mediante la aplicación de los conocimientos obtenidos y la evaluación de la propia aplicación de los mismos.

El método seguido durante el desarrollo se basa en el marco de trabajo *Agile* pensando en el desarrollo de la solución como un producto. De este modo se podrá iterar sobre una arquitectura definida y enfocarse en el desarrollo del producto en su totalidad, frente a un modelo tradicional de cascada donde cada fase de desarrollo recae en una sola parte del proyecto.

El formato de entrega será mediante un *minimum viable product* (MVP) en sprints de dos semanas durante el periodo de desarrollo y evaluación de las necesidades según la evolución del producto.

1.6. Planificación del Trabajo

Para la planificación del trabajo se han subdividido las tareas principales y se han mostrado en el siguiente diagrama de Gantt [1.1](#).

- Pec01 - Definición

- Pec01 - Planificación
- Pec02 - Búsqueda de fuentes del estado del arte
- Pec02 - Lectura de estado del arte
- Pec02 - Justificación de estado del arte
- Pec02 - Redacción
- Pec02 - Refinamiento de objetivos
- Pec03 - Planificación de Sprints
- Pec03 - Sprint 1
- Pec03 - Sprint 2
- Pec03 - Sprint 3
- Pec03 - Refinamiento/Redacción
- Pec04 - Revisión de apartados anteriores
- Pec04 - Redacción de nuevos apartados
- Pec05 - Presentación y defensa

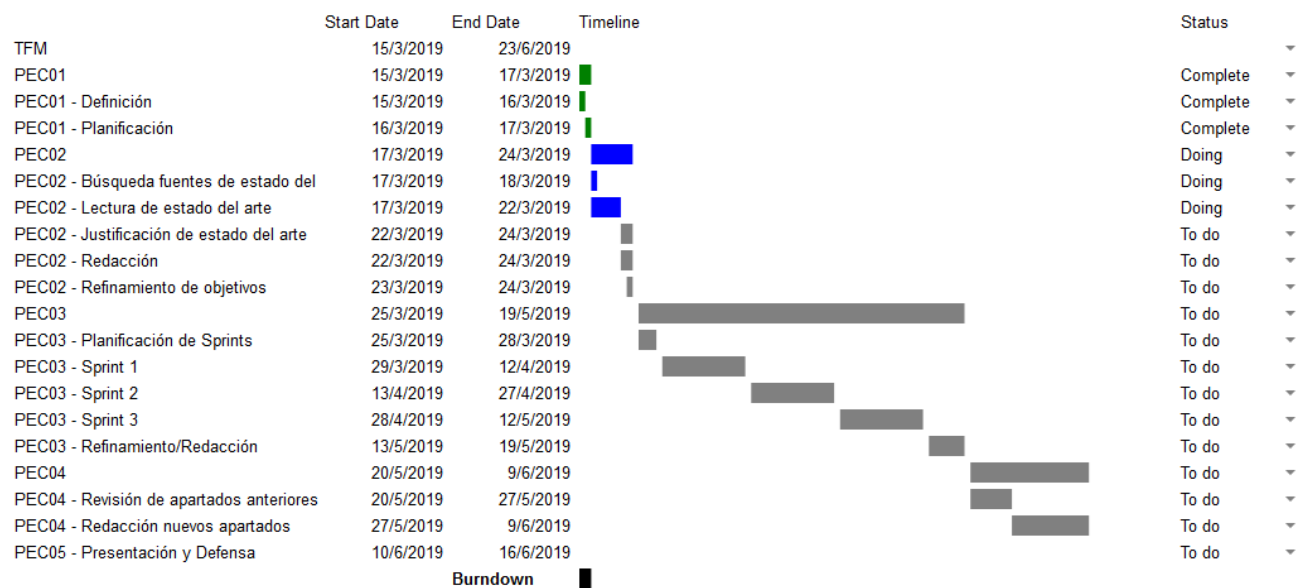


Figura 1.1: Planificación de tareas

Capítulo 2

Estado del Arte

Una de las aplicaciones más comunes dentro de la detección de anomalías es para la detección de intrusos (*Intrusion Detection*) y la creación de sistemas de detección de intrusos (*Intrusion Detection System, IDS*). Se trata de sistemas cuyo objetivo es la monitorización de los sistemas y redes informáticas, con el fin de alertar en caso de que puedan existir brechas de seguridad[1].

Con la ingente cantidad de redes y sistemas que se pueden encontrar a día de hoy es necesario incluir uno de estos sistemas, con el fin de mantener la integridad y disponibilidad de los mismos. Dentro de los IDS, se pueden identificar dos grandes implementaciones de estos sistemas, en los sistemas de detección de intrusos en Host (*Host-Based Intrusion Detection Systems, HIDS*) y los sistemas de detección de intrusos en red (*Network Intrusion Detection Systems, NIDS*).

- **HIDS:** Estos sistemas se caracterizan por implementarse en el Host utilizando información del propio sistema operativo para detectar actos maliciosos [2]. Esta información tiene distintos niveles de información, pero por lo general tienden a ser de bajo nivel sobre operaciones que se pueden estar realizando dentro del sistema. Esta información se consulta dentro de logs, por lo que el análisis de la información es más lento.
- **NIDS:** Para el segundo caso la monitorización se realiza a sistema de red, es decir, de comunicaciones entre distintos nodos y la monitorización de los paquetes que viajan entre ellos [1]. Esta información puede ser consumida en tiempo real, por lo que la reacción ante algún evento es más rápida que en los HIDS que necesitan revisar las acciones.

2.1. Métodos tradicionales de detección de anomalías

En la sección actual se describen algunos de los métodos tradicionales utilizados en IDS, tanto para HIDS como para NIDS.

- Network Security Monitor (NSM): se trata de uno de los primeros sistemas que permitió auditar el tráfico que circulaba dentro de la red [3]. El sistema escucha pasivamente dentro de la red y detecta si existe una conducta sospechosa al desviarse de patrones de conducta. La mayor parte de la monitorización se basa en protocolos estándar como *telnet*, *ftp*, *TCP/IP*, *etc.* por lo que le permitía utilizar una gran cantidad de datos heterogéneos.
- State transition analysis (USTAT): el sistema parte de que el host en un momento se encuentra en un estado seguro y que según las acciones que se realizan sobre el mismo el host cambia de estado, hasta que llega a un estado en el que compromete la seguridad [3]. Este sistema analiza los estados por los que ha pasado la máquina desde el estado seguro al comprometido.
- GrIDS: se trata de un IDS que utiliza un sistema de construcción de grafos basados en la red, donde cada nodo representa a un host y las aristas las conexiones entre los mismos. La representación gráfica de la actividad de la red permite ayudar al espectador en identificar qué está sucediendo [3].
- Haystack: en este caso el IDS se ayuda de métodos estadísticos para la detección de anomalías, definiendo estrategias para usuarios y grupos, además de definir variables del modelo como variables gaussianas independientes [4]. Para la detección se incluyen una serie de intervalos en los valores que en el momento que salen del rango normal, se calcula la distribución de probabilidades y si el *score* o puntuación es demasiado grande se genera una alerta.

Los métodos/sistemas listados se desarrollaron durante los años noventa, la tecnología ha evolucionado desde entonces y los sistemas se han vuelto más complejos y más propensos a los ciberataques. Por ello, se han desarrollado nuevas técnicas que se apoyan en el uso de técnicas de minería de datos (*Data Mining*) y las técnicas que se describirán a continuación de *Machine Learning*.

2.2. Machine Learning y la detección anomalías

El *Machine Learning* es una rama de la inteligencia artificial cuya premisa es hacer que la máquina aprenda una tarea sin haber sido específicamente programada para ello. El término fue descrito por Arthur L. Samuel en 1959 en un artículo en el que explica estudios de *Machine Learning* aplicado al juego de las damas [5], utilizando en una primera instancia métodos de aprendizajes más generales, como las redes neuronales de las que hablaremos más adelante, y otro métodos que tendrán que ser parametrizados para sus distintos usos.

La inteligencia artificial se puede entender como la inteligencia ejercida por las máquinas, al contrario que la inteligencia natural inherente a los humanos, que son capaces de realizar tareas cognitivas que permiten potenciar la resolución de las mismas e imitar comportamientos como el aprendizaje [6]. Dentro de la inteligencia artificial se pueden encontrar distintas definiciones según si se centran en el razonamiento o en el comportamiento:

- Actuar humanamente: este enfoque se basa en que las máquinas actúen como humanos más centrado en la interacción de máquinas con personas, más que en la resolución de problemas. Esta interacción se puede ver reflejada en el test ideado por Alan Turing en 1950, en el que se somete a la máquina inteligente.^a un interrogatorio realizado por un humano, donde éste no sabe que está hablando con una máquina, por lo tanto si es incapaz de detectar que se trata de una máquina esta ha actuado como un humano y puede considerarse inteligente.
- Pensar humanamente: esta vertiente se centra en imitar el pensamiento humano, entender cómo funcionan las mentes de los mismos y replicarlo en máquinas. Siguiendo esta corriente, si se consigue imitar el pensamiento humano, una máquina que resuelva el problema utilizará razonamientos humanos, en lugar de solucionar problemas a toda costa, independientemente de como lo realizan los humanos.
- Pensar racionalmente: basado en la lógica formal desarrollada en en siglos XIX y XX que permite formular los problemas en un lenguaje formal, utilizando el razonamiento matemático para resolver éstos.
- Actuar racionalmente: se centra en realizar el comportamiento más efectivo en un momento dado. Antes las distintas situaciones no siempre existe una acción correcta, pero si se puede llegar realizar una acción que minimice los riesgos.

Los enfoques mostrados conllevan sus propios enfoques filosóficos, sin embargo, han influenciado en gran medida a cómo se afrontan los problemas y cómo se han desarrollado las técnicas de inteligencia artificial que están en uso actualmente.

Como se ha comentado con anterioridad, el *Machine Learning* es una rama de la inteligencia artificial, dentro de la cual existe otra rama, *Deep Learning*. Las relaciones entre los términos se puede observar en la imagen 2.1 y como la inteligencia artificial engloba ambas ramas.

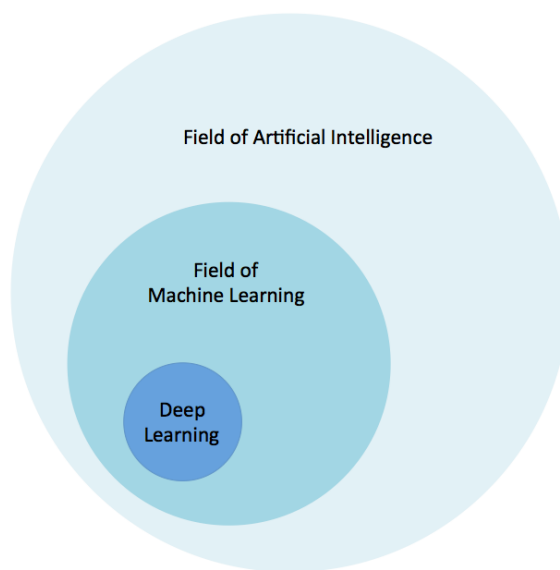


Figura 2.1: Relación entre Inteligencia Artificial, Machine Learning y Deep Learning

El *Deep Learning* se diferencia del *Machine Learning* convencional, en que están formados por modelos con distintas ramas de procesamiento que son capaces de extraer las representaciones de los datos con varios niveles de abstracción [7]. Mucha de la teoría relacionada se puede encontrar en los años 50 y 60, sin embargo, las aplicaciones y la capacidad de cómputo eran limitadas en la época, hasta ahora donde la capacidad de computación ha crecido a gran escala y en conjunto con las nuevas técnicas propuestas, una mayor cantidad de datos y la aplicación de transformaciones no lineales hacen que esta vertiente viva su época dorada.

Una de las implementaciones más representativas del *Deep Learning* se trata de las redes neuronales, las cuales se encuentran basadas en el funcionamiento del cerebro humano, utilizando el concepto de neuronas que pueden realizar cálculos, la conexión entre las mismas, la función de realizar una tarea específica, etc. Añadir, que se encuentra basado y que el cerebro humano es un sistema mucho más complejo que las redes neuronales y tiene muchos más comportamientos [8].

Dentro del *Machine Learning*, los algoritmos pueden dividirse según como se realice el entrenamiento del modelo, afectando también a la evaluación y las aplicaciones del mismo. Estas categorías son: *Aprendizaje Supervisado (Supervised Learning)*, *Aprendizaje No Supervisado (Unsupervised Learning)* y *Aprendizaje Semi-Supervisado (Semi-Supervised Learning)*

2.2.1. Aprendizaje Supervisado

El aprendizaje supervisado se caracteriza por utilizar un conjunto de variables de entrada y encontrar la función que más se aproxime a los valores de salida [9]. Actualmente es una de las metodologías más utilizadas en *Machine Learning* y es una de las más avanzadas en el campo, sin embargo, la necesidad de los valores de salida (variable dependiente o *target*) requiere que el conjunto de datos se encuentre etiquetado, es decir, para las observaciones de las variables de entrada debe existir la variable de salida para poder general el modelo. En muchos casos este etiquetado se debe de realizar manualmente, por lo que es un proceso que requiere mucho tiempo.

Por otro lado, cabe mencionar que el uso de métodos de aprendizaje supervisado tienen una mayor efectividad que el semi y no supervisado, sin embargo, la necesidad de tener las etiquetas y el coste que ello supone para grandes cantidades de datos, como es en el caso de la detección de anomalías, provocan que se busquen otras soluciones con aprendizajes semi y no supervisado. Además, suele suceder que en la detección de anomalías, éstas sean la menor parte de todos los datos, en otras palabras, la mayor parte de los datos son normales y solo una pequeña porción se trata de anomalías como tal, provocando que el conjunto de datos esté poco compensado.

Uno de los modelos más sencillos de aprendizaje supervisado es la regresión lineal (simple), ésta permite ajustar a una serie de puntos conocidos , x , y su respuesta y [10], generando una función que permite predecir los valores de y con valores no conocidos de x :

$$y = ax + b \tag{2.1}$$

En este caso se predicen los nuevos valores utilizando la función ajustada, existen distintos métodos para encontrar la recta que ajusta los puntos, siendo uno de los más utilizados el ajuste

por mínimos cuadrados:

$$a = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (2.2)$$

$$b = y - a\bar{x} \quad (2.3)$$

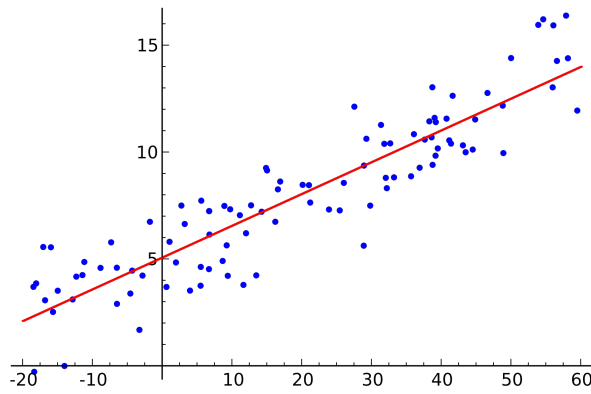


Figura 2.2: Regresión Lineal

Como se puede observar en la imagen 2.2, se ha generado una recta capaz de ajustarse a los puntos y predecir el valor para nuevos puntos. Para generar la función de la recta se ha necesita utilizar la variable dependiente y en el entrenamiento para poder generar la función que aproximará los nuevos puntos.

Durante los siguientes apartados se van a mostrar los distintos modelos utilizados en aprendizaje supervisado para la detección de anomalías, como las redes neuronales artificiales (*Artificial Neural Networks*, *ANN*, *Gradient Boosting Machines* y *Support Vector Machines*)

2.2.1.1. Feedforward Neural Network

Tal y como se ha comentado anteriormente, las redes neuronales intentan imitar el concepto de la conexión de neuronas del cerebro humano, éstas están formadas por una capa de entrada de datos (*Input*), una capa de salida (*Output*) y un conjunto de capas intermedias, también llamadas capas ocultas (*Hidden*) [8]. En las capas intermedias se encuentran las neuronas conectadas entre sí y pueden estar formadas desde una sola capa, hasta n capas, sin embargo, cuanto mayor sea

el número de capas, mayor coste de computación asociado. En la siguiente imagen 2.3 podemos ver como sería una arquitectura de una red neuronal:

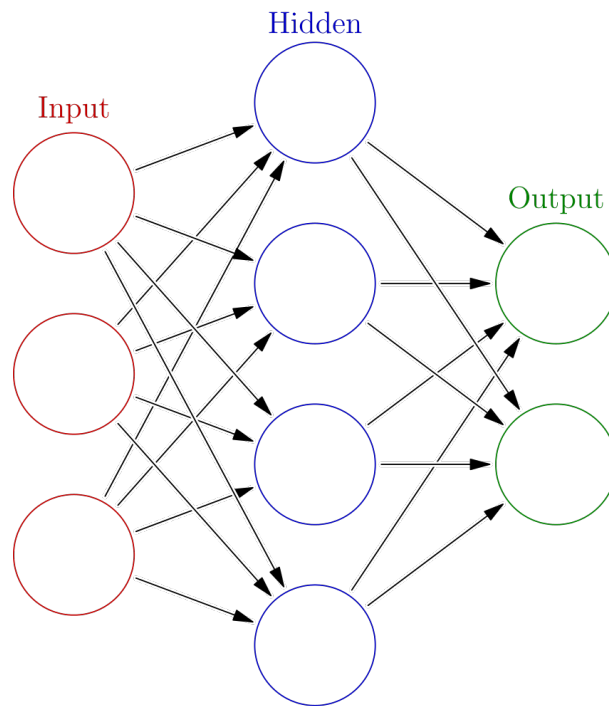


Figura 2.3: Red Neuronal Artificial

Dentro de las redes neuronales existe un caso base conocido como perceptron, un clasificador binario formado por la capa de entrada, una capa oculta de una sola neurona y una salida [11]. El perceptrón ayudará a incluir el concepto de los pesos (*weights*) y el *bias*.

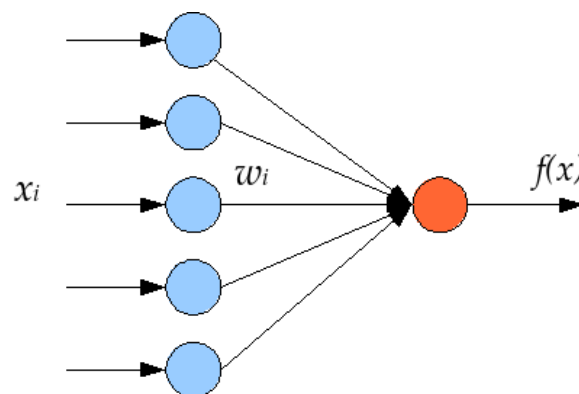


Figura 2.4: Arquitectura de un Perceptron

Los valores de entrada X son multiplicados por una serie de pesos W , en una primera instancia se inician aleatoriamente, y se realiza la suma de esta multiplicación, tras la cual se pasará a una función de activación. La función de activación define el valor de salida de la neurona y que se utilizará como valor de entrada para la siguiente neurona o como es en este caso como el valor de salida final. Un ejemplo de una función de activación es la función escalón, donde los valores menores de cero son iguales a cero y los mayores de cero son igual a uno. El *bias* se trata de un pequeño valor que modifica el output sin interactura con las neuronas. De este modo la clasificación de 0 o 1 en un perceptron sería así:

$$f(x) = g(\sum XW + b) \quad (2.4)$$

Donde $g(x)$ es la función de activación, X el vector de datos de entrada, W el vector de pesos y b el vector de bias.

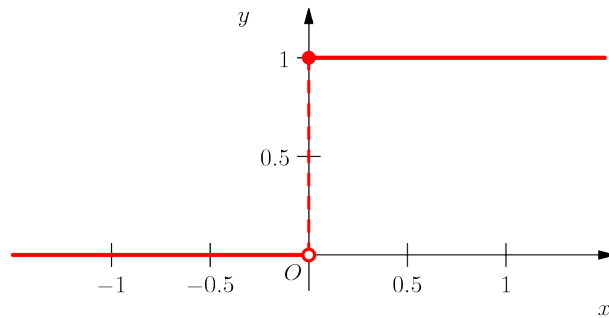


Figura 2.5: Función de activación: Función Escalón

Aumentando la complejidad de la arquitectura, es decir, el número de neuronas, el número de capas y distintas funciones de activación, se consiguen modelos más complejos pero mantienen el enfoque que se realiza en el perceptron.

Para optimizar la salida, se utilizan varias técnicas de optimización que permiten actualizar los pesos W y disminuir el error $E = f(x) - y$, una de las más utilizadas es el *backpropagation*, basado en el descenso del gradiente cuyo objetivo es obtener el gradiente de la función (vector de mayor pendiente) pero utilizar el valor opuesto del mismo multiplicado por un coeficiente de aprendizaje (*learning rate*), una analogía muy utilizada es imaginar como encontrar el punto más bajo del valle, para ello lo mejor es fijarse en que lugar del punto en el que nos encontremos tiene más pendiente negativa y seguir ese camino. El *backpropagation* utiliza este método, pero en lugar de calcular el gradiente sobre la función total (gran coste) se realiza el cálculo de los

gradientes locales, empezando desde el punto final de la función.

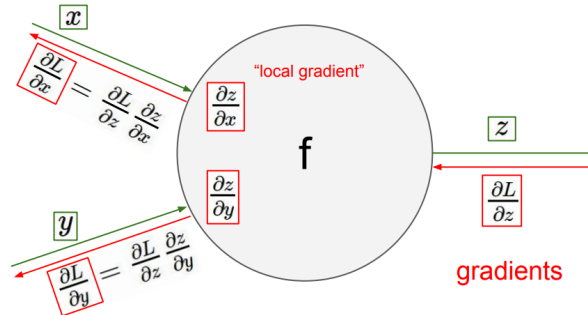


Figura 2.6: Backpropagation

La fórmula de la actualización de los pesos en una ANN sería similar a la siguiente:

$$W_{i+1} = W_i - \gamma dW \quad (2.5)$$

Donde γ es el *learning rate* y $L(x)$ es la función de pérdida.

La última capa, la capa de *outputs*, define el número de neuronas de salida, que puede ser igual al número de clases etiquetadas o una sola neurona en el caso de que existan solo dos clases, clasificación binaria. El valor de las salidas irá definido por la función de activación de la última capa de la capa oculta.

2.2.1.2. Gradient Boosting Machines

Las *Gradient Boosting Machines* son uno de los métodos más utilizados actualmente, dada la gran precisión que aportan a la resolución de los problemas. Estos modelos se basan en realizar el ensamblado de modelos débiles, normalmente árboles de decisión, para generar un predictor "fuerte" de una manera iterativa.

El término *Boosting* se refiere al método utilizado para mejorar el aprendizaje de un modelo mediante el ensamblado de modelos débiles, esto quiere decir que son modelos con una capacidad de predicción mejor que la pura aleatoriedad [12].

Concretamente el *Gradient Boosting* trata de generar un primer árbol para realizar las predicciones, es decir, crea un función $F(x)$ para aproximar y , posteriormente se utiliza la

función de coste definida (puede ser un error cuadrático medio para una regresión, por ejemplo) para ver como de bien ha realizado la predicción. Los residuales obtenidos se utilizan para crear un nuevo modelo que utilice estos residuales con el fin de volver a minimizar el error y finalmente se genera un nuevo modelo que utilice ambos. De esta forma en cada iteración se puede añadir n cantidad de nuevos modelos que ayuden a mejorar las predicciones mediante la corrección de los errores de los modelos previos [13].

Se crea el primer modelo

$$F_1(x) = y \quad (2.6)$$

Se crea un segundo modelo con los residuales

$$h_1(x) = y - F_1(x) \quad (2.7)$$

Se genera un nuevo modelo con los anteriores

$$F_2(x) = F_1(x) + h_1(x) \quad (2.8)$$

Dentro del ajuste a los residuales es donde entra la parte del gradiente, tal y como hemos explicado anteriormente se utiliza el método del descenso del gradiente (*Gradient Descent*) para la optimización de la función de pérdida. Para cada paso que se realiza, en vez de calcular los residuales como en la fórmula 2.7 se ajustan calculando el gradiente de la función de pérdida y ajustando un nuevo modelo, h , con los nuevos residuales y en la generación del nuevo modelo incluirle *gamma* (γ), similar al concepto explicado de *learning rate*:

$$r = -\left[\frac{\partial L(F(x_i), y)}{\partial F(x_i)}\right] \quad (2.9)$$

$$F_2(x) = F_1(x) + \gamma h_1(x) \quad (2.10)$$

En la práctica el uso de estos modelos se ha comprobado que es realmente eficaz para la resolución de problemas. Dentro de estos algoritmos se pueden destacar dos en concreto, cuyas implementaciones se han utilizado en varias soluciones dentro de las competiciones realizadas por Kaggle.

XGBoost

Diminutivo de *eXtreme Gradient Boosting* se trata de una nueva implementación de GBM realizada por Chen, Tianqi, et al [14], que demuestra una mayor escalabilidad y velocidad que los métodos tradicionales, utilizando una menor cantidad de recursos de los sistemas.

LightGBM

Desarrollado por Microsoft [15] se propone mejorar los problemas encontrados en algoritmos como *XGBoost* cuando existe una alta cantidad de variables, esto es ocasionado por la necesidad de evaluar el punto óptimo para la división de los árboles de decisión. Propone su propia forma de efectuar la ganancia de información (por división) y demuestra una mejora en el tiempo de cómputo sobre los modelos convencionales.

Algunas de las competencias donde se han utilizado los algoritmos (también ensamblándolos) son las siguientes:

- Home Credit Default Risk: Can you predict how capable each applicant is of repaying a loan? [16].
- Corporación Favorita Grocery Sales Forecasting: Can you accurately predict sales for a large grocery chain? [17].
- Google Analytics Customer Revenue Prediction: Predict how much GStore customers will spend [18].

2.2.1.3. Support Vector Machines

Se trata de otro algoritmo enfocado a la clasificación mediante la separación de los puntos por un hiperplano [10]. Por ejemplo, para un set de datos etiquetados con dos categorías (clasificación binaria) se busca encontrar la línea (hiperplano) que mejor separe ambos puntos, como puede observarse en la imagen 2.7 existen distintos hiperplanos (H_1, H_2, H_3) que separan las distintas clases, como pueden observarse tanto H_1 como H_2 son capaces de separar las dos clases por completo, mientras que H_3 fallaría en la tarea de clasificación.

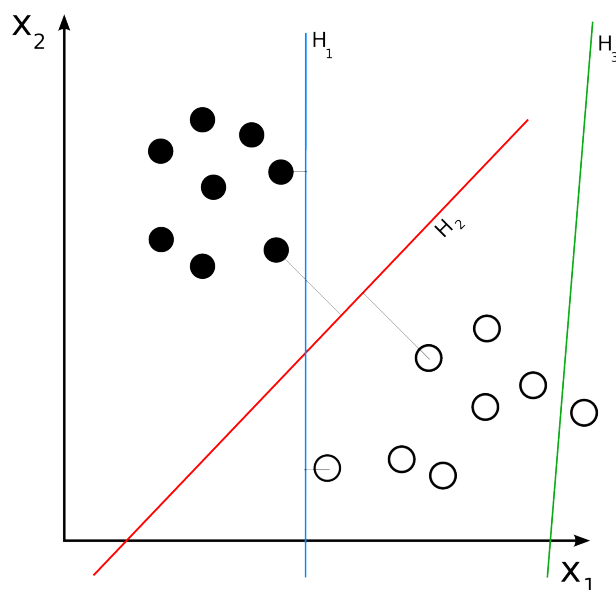
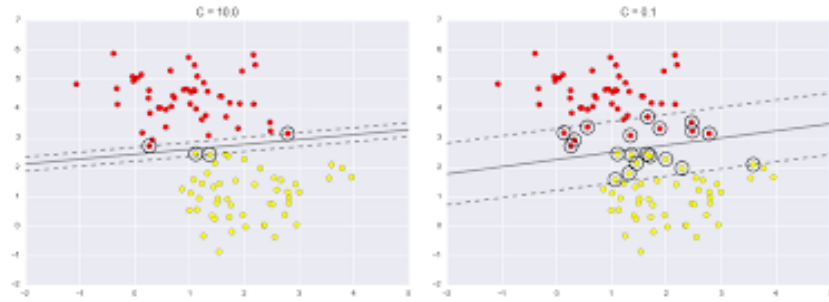


Figura 2.7: Hiperplanos SVM

Se puede definir el hiperplano como la "línea" que puede dividir los datos entre las dos clases, es decir, se pretende separar el espacio en dos mitades. Esta "línea" para un problema de p dimensiones sería de $p - 1$ dimensiones, en otras palabras, para un problema de dos dimensiones sería una línea, para uno de tres sería un plano y así sucesivamente hasta n dimensiones. Esta separación indica que según donde se encuentren los puntos con respecto al hiperplano se clasificarán como una clase u otra.

$$\beta_0 + X_1\beta_1 + X_2\beta_2 + \dots + X_n\beta_n = 0 \quad (2.11)$$

Como se ha mencionado en la imagen 2.7, existen varios hiperplanos, sin embargo, se debe de buscar el hiperplano óptimo de los posibles, siendo este el que se encuentra más alejado de los puntos de entrenamiento, es decir, es el hiperplano con mayor margen entre los puntos y el hiperplano. Sin embargo, si solo nos basamos en encontrar el mejor hiperplano según los márgenes se puede dar el caso de que un nuevo punto afecte en gran medida al hiperplano, por maximizar los márgenes, por ello se utiliza un parámetro C que permite "ablandar" el clasificador y permitir que ciertos puntos puedan encontrarse entre el hiperplano y el margen, la mala clasificación de unos pocos puntos permitidos evita sobreentrenar el modelo y permitir generalizar mejor.

Figura 2.8: Hiperplanos para distintos valores de C

Por ahora, se ha hablado de separar los datos linealmente y en dos clases, sin embargo, ese no es siempre el caso, podemos tener más de dos clases o por ejemplo tenemos datos que no sean separables por una "línea", pensemos en dos conjuntos de puntos que forman una serie de circunferencias concéntricas (2.9), para este caso una línea nunca logrará separar bien los puntos. Para esta última, la *SVM* puede utilizar los *kernels*, se trata de funciones que permiten cambiar la dimensionalidad del espacio para ayudar a la separación de las clases.

Kernel lineal

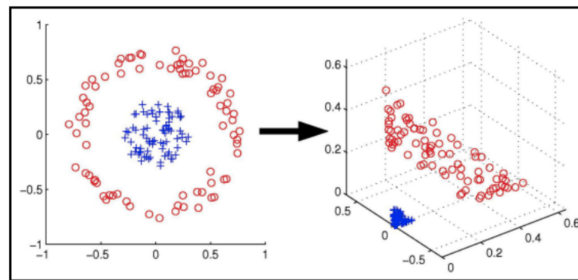
$$K(x_i, x_{i'}) = \sum_{j=1}^p x_{ij} x_{i'j} \quad (2.12)$$

Kernel polinomial

$$K(x_i, x_{i'}) = \left(1 + \sum_{j=1}^p x_{ij} x_{i'j}\right)^d \quad (2.13)$$

Kernel radial

$$K(x_i, x_{i'}) = \exp\left(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2\right) \quad (2.14)$$

Figura 2.9: Efecto de aplicación de *Kernel* radial en SVM

Para el caso de tener más de dos clases se pueden aplicar dos enfoques distintos:

- Clasificación uno contra uno: para este caso se crean una SVM por cada par de clases y realiza la clasificación para ellos, luego el resultado es elegido observando cuantas veces el punto cae en una categoría para las distintas SVM clasificándolo a la clase más frecuente.
- Clasificación uno contra todos: en este caso se realiza una SVM por clase, de modo que se compare la clase escogida con el resto y valorando en cual de los modelos generados es más apta la clasificación del punto.

2.2.2. Aprendizaje No Supervisado

Este nuevo tipo de aprendizaje se distingue porque no necesita ese conjunto de etiquetas que se utiliza por observación, es decir, en este caso no se tiene el valor de la clase o el valor numérico correspondiente a cierto conjunto de datos. Para este caso no se enfoca a la predicción de los valores, debido a que no tenemos esa variable *target* para la que queremos desarrollar la función que la aproxime lo máximo posible, en este caso, el objetivo principal es intentar descubrir nuevos patrones o información oculta dentro de los datos [10].

Muchas veces el uso del aprendizaje no supervisado no se encuentra con el simple objetivo de predecir, si no que puede ser parte de un proceso de análisis de los datos, de modo que permitan encontrar más información útil de los mismos, como por ejemplo, se puede realizar una segmentación de usuarios, de modo que posteriormente se pueda estudiar que causa esta segmentación de los mismos y mejorar los servicios que se les puedan ofrecer.

El uso del aprendizaje no supervisado esta creciendo cada vez más, dado que en la mayor parte de los casos los datos no vienen etiquetados, tal y como se ha comentado anteriormente el etiquetado de los mismos es un proceso costoso y además mientras que en el supervisado se puede "supervisar" si el resultado es correcto en la fase de entrenamiento, mientras que para el no supervisado, no se tiene esta capacidad de conocer cual es el valor real, dado que es imposible saber cual es la respuesta correcta.

2.2.2.1. Análisis de Componentes Principales

Se trata de uno de los algoritmos más conocidos cuya finalidad es explicar en una serie de "componentes", siendo la cantidad de éstos menor a la cantidad de variables de los datos iniciales, que permiten explicar la mayor parte de la varianza de los datos originales. Estos "componentes" se trata de un conjunto de nuevas variables que son capaces de explicar una gran parte de la varianza de un conjunto mayor, en otras palabras, permite obtener gran parte de la varianza explicada en una dimensionalidad menor.

En las siguientes imágenes se ha utilizado un set de datos sintéticos de dos dimensiones, en el que se muestra como en la imagen de la izquierda 2.10 se encuentran nubes de datos que cambian en función de los ejes x e y y como la aplicación del análisis de componentes principales 2.11, *PCA* de ahora en adelante, genera nuevas variables (dos componentes principales) que tienden a explicar la varianza de las variables originales. En este caso se puede apreciar como la aplicación del algoritmo "mueve" las observaciones, haciendo que así sean más identificables en a lo largo de los ejes.

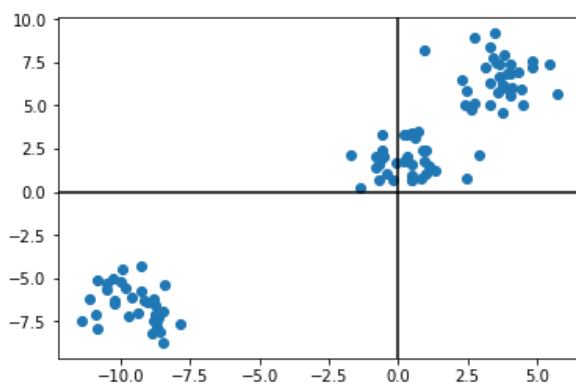


Figura 2.10: Datos de dos variables sin PCA

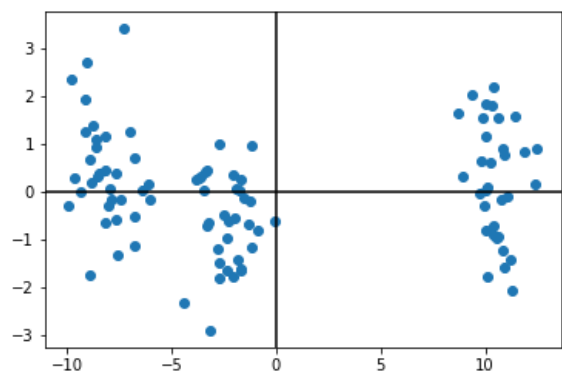


Figura 2.11: PCA aplicado a dos variables

Este ejemplo se trata de un caso sencillo, por lo general esta técnica también se puede utilizar para visualizar los datos de una gran dimensionalidad en dos o tres variables, utilizando ese número de componentes principales, que se tratan de una combinación lineal de las variables anteriores [10].

Para realizar el cálculo del PCA, se puede realizar mediante el uso de álgebra lineal, más en concreto con una descomposición de autovalores y autovectores que siguen los siguientes pasos

para una matriz X con p variables y n observaciones [19]:

Se calcula la media de cada una de las columnas

$$\bar{X} = \bar{X}_1, \bar{X}_2, \dots, \bar{X}_p \quad (2.15)$$

Se normaliza el valor de las columnas restando la media de cada columna a todos los valores de la matriz

$$X_{norm} = X - \bar{X} \quad (2.16)$$

Se genera la matriz de covarianza con la matriz normalizada, siendo E la media, por cada par de variables [20]

$$Cov[X_i, X_j] = E[(X_i - E[X_i])(X_j - E[X_j])] \quad (2.17)$$

De la matriz de covarianza se calculan los autovalores y autovectores

$$\lambda, \nu < -X_{cov} \quad (2.18)$$

Una vez obtenidos los autovectores y autovalores, se seleccionarán una cantidad m que será la cantidad de componentes principales elegidos

$$B = [\lambda_1 \dots \lambda_m, \nu_1 \dots \nu_m] \quad (2.19)$$

Y con estos calculamos la proyección P de la matriz original X con m componentes principales

$$P = B^T \cdot A \quad (2.20)$$

2.2.2.2. Agrupación

Tal y como el propio nombre indica, actualmente se utilizan métodos de agrupación en el que podemos encontrar grupos o *clusters*. Esta agrupación se realiza mediante la similitud entre observaciones, donde la similitud puede ser la distancia euclídea entre las observaciones [10], útil cuando todas las variables son numéricas. En la imagen 2.12 podemos ver un ejemplo donde datos (sin etiquetar, imagen de la izquierda) se analizan y se agrupan (imagen de la derecha). En este caso la agrupación es clara, sin embargo, la agrupación se puede realizar con datos más abstractos y menos representables dada su alta dimensionalidad.

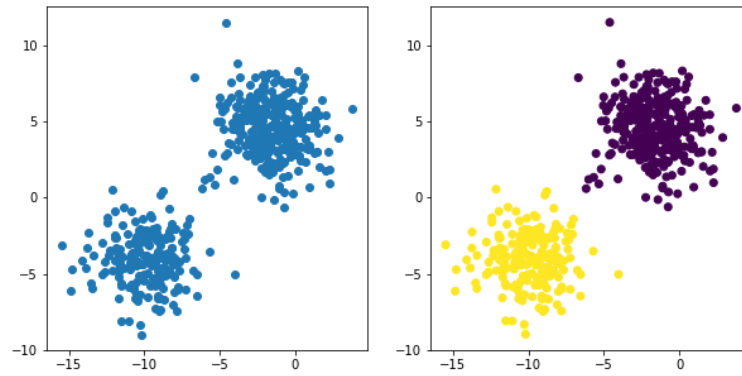


Figura 2.12: Ejemplo de agrupamiento

Como se ha comentado, la similitud entre observaciones ayuda a encontrar que grupos son más similares, por lo tanto dependiendo del problema se pueden utilizar distintas distancias como la euclídea o manhattan, siendo estas de las más utilizadas:

Distancia euclídea

$$d(X_1, X_2) = \sqrt{(X_{11} - X_{21})^2 + (X_{12} - X_{22})^2 + \dots + (X_{1n} - X_{2n})^2} \quad (2.21)$$

Distancia manhattan

$$d(X_1, X_2) = \sum_{i=1}^n |X_{1i} - X_{2i}| \quad (2.22)$$

Estas distancias se aplican a datos con variables numéricas por lo que no son adecuadas cuando existe variables discretas, sin embargo, se pueden realizar transformaciones en estas

variables a numéricas para poder utilizarlas. En el caso donde la mayor parte de las variables sean categóricas esta transformación no daría buenos resultados, por lo que se pueden utilizar distancias específicas para variables discretas, como *Goodall*, *Lin* o *Smirnov* [21].

Dentro de los algoritmos de clusterización los más conocidos son los siguientes:

- **KMeans o K-medias:** es un algoritmo de clusterización en el cual se indica a priori el número de clústers K que se quieren realizar [10]. Una vez definido se generan k centroides aleatorios, se calculan las distancias de estos centroides al resto de puntos de modo que se asocia a cada punto el centroide que más cercano se encuentre (se puede decir que se le asigna una clase o grupo) y finalmente se actualiza el centroide utilizando los valores medios de todos sus puntos asociados. Se vuelven a calcular las distancias, asignar de nuevo los centroides y la actualización de éstos hasta que no se realizan grandes cambios, es decir, converge o hasta un número máximo de rondas.
- **DBSCAN:** este algoritmo se basa en la generación de los grupos mediante la identificación de áreas de alta densidad (de puntos) frente a zonas de poca densidad [22]. En este caso no se necesita indicar el número de clústers a priori, pero sí un número mínimo de observaciones como para poder considerar que se trata de un nuevo grupo. En este algoritmo no todas las observaciones se agrupan, si no cumplen la cantidad mínima estas observaciones se pueden considerar como *outliers* o anomalías.
- **Algoritmos jerárquicos:** en este tipo de algoritmos se define una medida de disimilitud (distancia euclídea, por ejemplo) entre los distintos puntos (considerados como un grupo cada uno), de modo que aquellos puntos con la menor disimilitud se agrupan en un nuevo clúster, generando un nuevo grupo de puntos. Con este grupo de puntos surge una nueva medida que permite identificar la disimilitud entre los mismos, *linkage*. Por lo general existen distintos modos de *linkage* siendo los más comunes: completo (se calcula la disimilitud entre todos los puntos de los grupos y se selecciona la mayor), *single* (se selecciona la menor), media (se hace la media de todas las disimilitudes) y centroide (se calcula la disimilitud entre los centroides de los grupos)

Tal y como se ha comentado anteriormente, generalmente se suelen utilizar conjuntos de varias técnicas que permitan mejorar los resultados obtenidos utilizando una técnica sola, un ejemplo en la detección de anomalías es el uso de PCA y KMeans, arrojando buenos resultados en dispositivos IOT [23].

2.2.2.3. Autoencoders

Los autoencoders son redes neuronales con una determinada estructura, de modo que no se necesite un set de datos etiquetados (variable dependiente o *target*), el objetivo de esta red neuronal es reconstruir los datos de entrada en la capa de salida. Esto se realiza utilizando el mismo número de neuronas de entrada en la salida, de modo que utilizando las técnicas de las redes neuronales, como la optimización por descenso del gradiente y el *backpropagation* se puedan generar nuevas variables (*features*) que reconstruyen los datos de entrada [24].

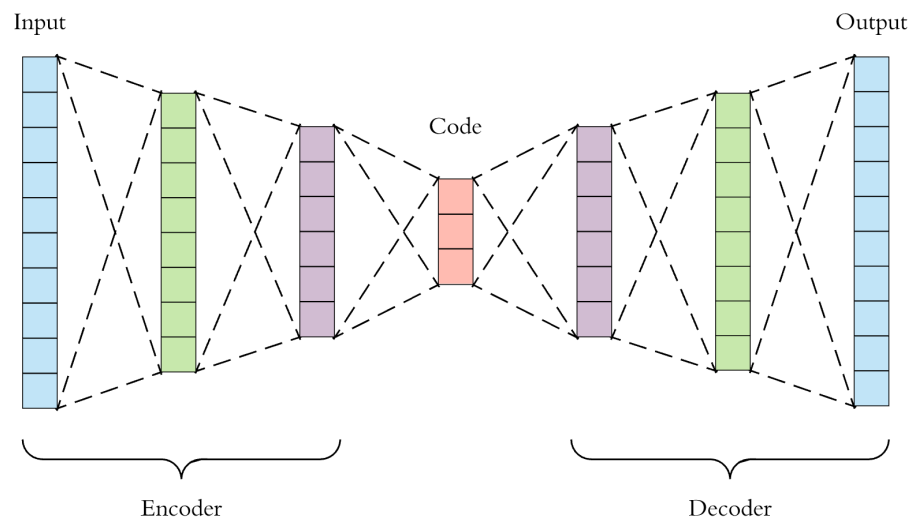


Figura 2.13: Arquitectura de un Autoencoder

Como se puede observar en la imagen 2.13 las dimensiones de entrada son iguales en la salida, y las capas intermedias se pueden dividir en dos partes:

- *Encoder*: permite codificar la entrada a una dimensionalidad menor, de modo que se genera una representación de menor dimensionalidad de la entrada y se generan las nuevas *features*. La capa intermedia se puede utilizar como representación de menor dimensionalidad, similar a la reducción de dimensionalidad del PCA, que puede utilizarse como variables para otro algoritmo.
- *Decoder*: decodificador utilizando la capa intermedia donde se pretende reconstruir los datos de entrada, utilizando las *features* generadas por el *Encoder* hasta llegar a la dimensionalidad de entrada.

Utilizando este método es posible, a partir de datos sin etiquetar, entrenar una red neuronal que permita reconstruir fielmente los datos de entrada, pero en el caso de las anomalías esta reconstrucción no sería del todo fiel (o similar a datos comunes), en otras palabras, el error de la reconstrucción sería mayor que el de una representación normal. De este modo se permite identificar picos en el error que muestren que la reconstrucción no es fiel y que la observación en cuestión se pueda tratar de una anomalía.

En la imagen 2.14 se puede observar un ejemplo del uso de autoencoders que permite clasificar una transacción fraudulenta (clase 1) de una no fraudulenta (clase 0). En este caso la reconstrucción de los datos no fraudulentos representan un error muy pequeño, en comparación a sus contrapartidas, y se puede observar como la mayor parte de estos se encuentran debajo del umbral definido. Aquellos datos anómalos no acaban de reconstruirse con total fidelidad aumentando el error definido (en este caso el error cuadrático medio). No todos son identificados correctamente, pero si permite aumentar la detección de los mismos, en algunos casos donde la cantidad de datos es muy grande y la cantidad de anomalías es muy pequeña este tipo de enfoques puede ayudar a disminuir la cantidad de trabajo necesario para confirmar el fraude, un ejemplo sería en el blanqueo de capitales donde la mayor parte de casos no constituyen blanqueo, reducir el número de casos a revisar puede ahorrar una gran cantidad de costes y tiempo.

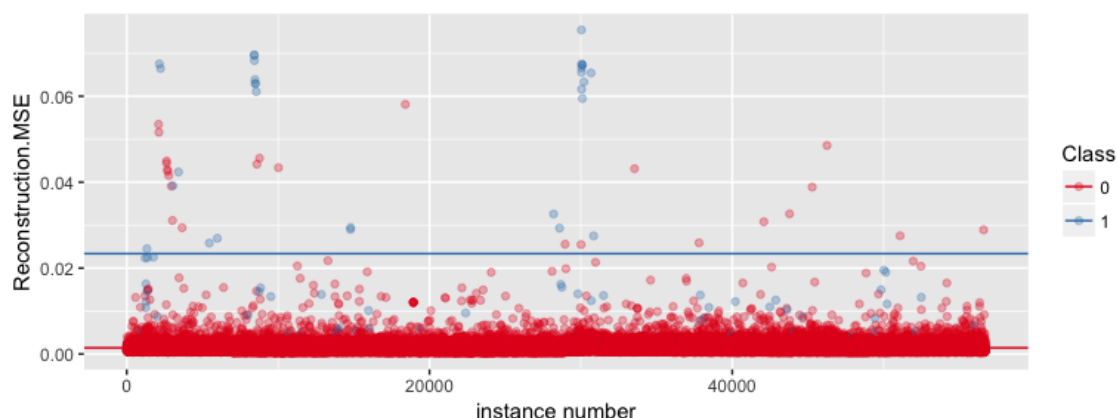


Figura 2.14: Ejemplo de detección de anomalías con Autoencoders

Capítulo 3

Desarrollo

En esta parte de la memoria se detallarán los pasos seguidos durante el desarrollo y la aplicación de los algoritmos de *Machine Learning* sobre los datos utilizados. Asimismo, se profundizará en los datos utilizados y cómo se han preprocesado para su posterior aplicación en los algoritmos.

Generalmente, la mayor parte del tiempo de un proyecto de *Data Science* se invierte en la extracción y preprocesado de los datos, hasta que tengan la suficiente calidad como para arrojar resultados fiables tras aplicar los algoritmos. Este ha sido el caso del presente proyecto, la mayor parte de los esfuerzos han recaído en la comprensión de los datos y en cómo hacer frente (tanto por tiempo como por tamaño) a los datos disponibles.

Tanto el preprocesado como la aplicación de los algoritmos de *Machine Learning* se han realizado utilizando el lenguaje de programación *Python* en su versión 3.6 y la exploración y análisis de los datos se han realizado en *Jupyter Notebooks* para facilitar la visualización de los datos y las acciones tomadas. Todo estos procesos se han realizado en una máquina remota proporcionada por la UOC con 10 núcleos y 64gb de memoria RAM.

3.1. Datos utilizados

Para este proyecto se va a utilizar un *dataset* con información de las conexiones entre dispositivos IoT, mediante el análisis de los paquetes obtenidos por un analizador de paquetes o *sniffer*, un software que permite observar y recopilar las acciones que desarrollan dentro de una red [25]. Utilizando los datos recopilados por este software se preprocesarán de modo que se pueda extraer la información y posteriormente transformarla y crear nuevas variables.

Los datos iniciales fueron preprocesados inicialmente con la información de las anomalías (únicamente anomalías) que se habían recopilado en un espacio de tiempo de dos meses. En este caso los datos solo contaban con tráfico ilegítimo, por lo que no se podía hacer una detección de anomalías como tal.

Posteriormente, se proporcionaron los datos en bruto en un fichero de extensión *pcap* resultante del *sniffer* mencionado anteriormente. Este fichero contiene información de las distintas capas en las que se almacenan los datos, siendo las de más interés las capas *TCP* e *IP*.

- *Transmission Control Protocol*: originario de los años 70/80, consiste en un protocolo para la transmisión de paquetes en conexiones *host* a *host* en redes [26].
- *Internet Protocol*: protocolo que permite la transmisión de datos de origen a destino donde el origen y el destino son direcciones de longitud fija y proporciona los medios para la división y reagrupación de grandes datos en grupos más pequeños [27] .

Sobre estas capas se obtiene la información disponible al más bajo nivel y posteriormente se generarán nuevas variables en base a los datos obtenidos. Siendo los datos de la capa *TCP* los siguientes:

- Puerto de origen: indica el puerto de origen utilizado para realizar la conexión desde la IP de origen.
- Puerto de destino: indica el puerto de destino utilizado en la conexión entre *hosts*.

Los datos obtenidos de la capa IP:

- IP de origen: dirección con longitud fija del origen de la conexión.
- IP de destino: dirección con longitud fija del destino de la conexión.
- *Length*: longitud/tamaño de los paquetes enviados durante la conexión.

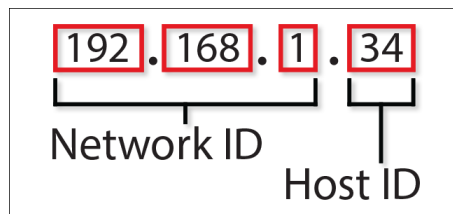


Figura 3.1: Ejemplo de dirección IPv4

Otros datos:

- *Payload*: contenido de los datos transmitidos, pudiendo ser texto, imágenes o comandos a ejecutar en el *host* de destino, con posibilidad de estar cifrado.
- *Timestamp*: marca de tiempo que indica en qué momento se realizó la conexión.

Un ejemplo más visual de los datos a bajo nivel se puede realizar utilizando la herramienta *Open Source, Wireshark*, fundamental a la hora de entender cómo están distribuidos los datos en las capas.

No.	Time	Source	Destination	Protocol	Length	Info
1	0.000000	104.248.135.165	133.34.156.76	TCP	66	47307 → 23 [ACK] Seq=1 Ack=1 Win=1173 Len=0 TSval=216133704 TSecr=4294855037
2	0.000415	133.34.156.67	104.248.135.165	TELNET	1514	Telnet Data ...
3	0.000663	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
4	0.000671	104.248.135.165	133.34.156.67	TCP	66	50293 → 23 [ACK] Seq=1 Ack=4294939225 Win=846 Len=0 TSval=216133697 TSecr=4294940066
5	0.000933	104.248.135.165	133.34.156.76	TCP	66	47307 → 23 [ACK] Seq=1 Ack=2897 Win=1173 Len=0 TSval=216133705 TSecr=4294855038
6	0.001681	104.248.135.165	133.34.156.76	TCP	66	47307 → 23 [ACK] Seq=1 Ack=7241 Win=1173 Len=0 TSval=216133706 TSecr=4294855038
7	0.002266	104.248.135.165	133.34.156.67	TCP	66	50293 → 23 [ACK] Seq=1 Ack=4294942121 Win=846 Len=0 TSval=216133699 TSecr=4294940067
8	0.002434	133.34.156.67	104.248.135.165	TELNET	1514	Telnet Data ...
9	0.002485	104.248.135.165	133.34.156.67	TCP	66	50293 → 23 [ACK] Seq=1 Ack=4294945017 Win=846 Len=0 TSval=216133699 TSecr=4294940067
10	0.002624	133.34.156.67	104.248.135.165	TELNET	1514	Telnet Data ...
11	0.002706	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
12	0.002877	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
13	0.003069	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
14	0.003262	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
15	0.003453	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
16	0.003959	104.248.135.165	133.34.156.76	TCP	66	47307 → 23 [ACK] Seq=1 Ack=10137 Win=1173 Len=0 TSval=216133708 TSecr=4294855040
17	0.004050	133.34.156.76	104.248.135.165	TELNET	1514	Telnet Data ...
Protocol: TCP (6) Header checksum: 0x3f51 [validation disabled] [Header checksum status: Unverified] Source: 133.34.156.67 Destination: 104.248.135.165 Transmission Control Protocol, Src Port: 23, Dst Port: 50293, Seq: 1, Ack: 1, Len: 1448 Source Port: 23 Destination Port: 50293 [Stream Index: 1] [TCP Segment Len: 1448] Sequence number: 1 (relative sequence number) [Next sequence number: 1449 (relative sequence number)] Acknowledgment number: 1 (relative ack number) 1000 ... = Header Length: 32 bytes (8) > Flags: 0x010 (ACK) Window size value: 5792 [Calculated window size: 5792] [Window size scaling factor: -1 (unknown)]						
0000	00 10 00 00	00 00 00 00	00 00 00 00	00 00 00 00	00 00 00 00	00 00 00 00
0010	05 dc a5 c7	40 00 3e 06	3f 51 85 22	9c 43 68 f8
0020	87 a5 00 17	c4 75 b2 76	19 7d fd cc	c7 b6 80 10
0030	16 a0 89 34	00 00 01 01	00 0a ff ff	95 e6 0c e1
0040	f0 3e 1c 3c	a4 25 9c 27	21 e0 99 03	d8 ff bd 27
0050	24 00 bf af	20 00 02 af	1c 00 b1 af	18 00 b0 af
0060	10 00 bc af	21 90 00 00	10 00 00 10	01 00 11 24
0070	0e 00 20 12	00 00 00 00	50 8b 99 8f	00 00 00 00
0080	20 f8 20 03	21 20 00 02	10 00 bc 8f	07 00 40 14
0090	79 00 02 24	04 00 02 12	ff ff 31 26	59 00 02 24
00a0	02 00 02 16	00 00 00 00	01 00 52 26	fc 88 82 8f
00b0	00 00 00 00	00 00 44 8c	00 00 00 00	34 00 82 8c
00c0	00 00 00 00	0d 00 40 10	00 00 00 00	10 00 83 8c
00d0	18 00 82 8c	00 00 00 00	2b 10 62 00	04 00 40 10
00e0	01 00 62 24	00 00 70 90	0d 0a 00 00	10 10 00 82
00f0	ac ac 8a 99	8f 03 00 00	10 00 00 00	0c 05 99
0100	8f 00 00 00	00 20 20 20	20 20 20 20	f8 20 03 00
0110	00 00 10 00	bc 8f 21 00	40 00 ff ff	02 24 03 00
0120	02 12 0d 0a	00 02 24 d3	ff 02 16 00	00 00 21
0130	10 40 02 24	00 bf 8f 20	00 b2 8f 1c	00 b1 8f 18

Figura 3.2: Ejemplo de datos pcap en Wireshark

El *Timestamp* puede indicar que el problema debería tratarse como una serie temporal y generar variables en base a esta característica.

Los datos provienen de la recolección de información realizada el 2 de Febrero de 2019, contando tanto con tráfico legítimo como ilegítimo. Sin embargo, no existe clasificación de los datos, por lo que se trata de un problema de carácter no supervisado y como consecuencia se utilizarán técnicas como las descritas en el capítulo 2 de aprendizaje no supervisado.

3.2. Exploración y Preprocesamiento

En este apartado se van a mostrar los pasos seguidos durante la exploración y el procesamiento así como la información obtenida durante el proceso. Ambas etapas se realizan de manera conjunta, es decir, la exploración conlleva la realización de distintos preprocesamientos o otras palabras, se realiza *Feature Engineering*. Se puede definir como el proceso de exploración y procesamiento de datos que ayudan a la creación de nuevas variables que puedan mejorar la eficiencia del modelo[28].

El paso inicial era la extracción de los datos del fichero *pcap*, el tamaño total era cercano a los 14 GB, de modo que en el momento de cargar la información en memoria, el uso aumentaba considerablemente dificultando su manipulación, por lo que el primer paso fue investigar si existía alguna forma de dividirlo. El resultado fue el uso de la herramienta *tcpdump*, una línea de comandos utilizada para el análisis de paquetes escrita en C/C++ [29], la misma permitía la división del fichero original en pequeños ficheros con el tamaño deseado. La división se realizó con un tamaño máximo de 500 MB, para permitir la posterior paralelización de los ficheros y evitar un uso muy alto de memoria, dando como resultado 15 ficheros.

Una vez obtenidos los ficheros se deben procesar de modo que el resultado sean datos de forma tabular, que puedan ser explorados y preprocesados con mayor facilidad. Para realizar este proceso se utilizó el módulo de *python*, *scapy*, utilizado comúnmente en tareas de análisis de redes. Gracias a la granularidad de los nuevos ficheros y a utilizar *python* se procesan 9 ficheros simultáneamente utilizando el módulo mencionado, éste se encarga de recorrer las capas mencionadas y extraer los datos y organizarlos de forma tabular, concretamente en *csv*. En ciertos casos, los paquetes pueden contener el *payload*, cuando este es el caso se extrae y se añade como una columna más, para el resto se trata como si fueran *NA* o datos sin registro.

Una vez terminado el proceso se obtienen los primeros datos en brutos de carácter estructurado. Se genera una cantidad análoga de ficheros con un volumen tres veces superior a los datos originales, con las variables mencionadas en la sección anterior. Antes de realizar ningún tipo de exploración se realiza un preprocesado en paralelo para generar nuevas variables sencillas que puedan aportar información de cara a la primera exploración. Las nuevas variables son las siguientes:

- Red 1 IP origen

- Red 2 IP origen
- Red 3 IP origen
- Host IP origen
- Red 1 IP destino
- Red 2 IP destino
- Red 3 IP destino
- Host IP destino
- is_busybox
- is_sh
- is_enable

Las variables en referencia a las IPs se generaron haciendo una división de las direcciones:

$$192,168,1,1 \Rightarrow 192_{red1} \ 168_{red2} \ 1_{red3} \ 1_{host} \quad (3.1)$$

El resto de las variables se obtuvieron mediante el *parseo* del *payload* en aquellos registros que así lo tuvieran y buscando la existencia de las palabras clave *sh*, *enable* y *busybox*, siendo la mayor parte de los registros sin contenido en el *payload* o incluso cifrado. La razón por la que se eligieron estas palabras clave es por el uso de este tipo de comandos en los dispositivos IoT por *malwares* conocidos, *busybox* se trata de una versión mínima de un sistema operativo *UNIX* utilizado en pequeños dispositivos y que puede ser explotado por *malware*[30]. De la misma manera *sh* y *enable* son comandos comúnmente utilizados en este tipo de ataques. Las tres variables son de tipo binario, es decir, indica si en el registro existe el comando o no.

El proceso de exploración se realiza en *Jupyter Notebooks* con la ayuda del módulo *matplotlib* y *seaborn* para las visualizaciones. Para esta primera exploración, se van a buscar los medios más utilizados en las distintas conexiones, es decir, las IPs más utilizadas, puertos más utilizados, información de las variables binarias y la evolución de la cantidad de conexiones en el tiempo.

Timestamp	IP1	IP1_0	IP1_1	IP1_2	IP1_3	IP2	IP2_0	IP2_1	IP2_2	IP2_3	Length	Payload	Port1	Port2	is_busybox	i
2019-02-01 09:37:56.421462	71.6.146.185	71	6	146	185	133.34.156.185	133	34	156	185	44	b'\x00\x00'	24858	8649	0	
2019-02-01 09:37:56.446036	185.176.27.114	185	176	27	114	133.34.156.140	133	34	156	140	40	b'\x00\x00\x00\x00'	52462	6334	0	
2019-02-01 09:37:56.512569	104.248.37.237	104	248	37	237	133.34.156.78	133	34	156	78	52	NaN	43082	23	0	
2019-02-01 09:37:56.549846	209.141.57.185	209	141	57	185	133.34.156.140	133	34	156	140	60	NaN	32860	2323	0	
2019-02-01 09:37:56.550018	151.73.61.209	151	73	61	209	133.34.156.65	133	34	156	65	40	b'\x00\x00\x00\x00'	45431	23	0	

Figura 3.3: Datos en forma tabular

Para realizar la exploración se unen todos los datos contenidos en los ficheros generados ordenados por la marca de tiempo. Una vez cargados, se realizan las operaciones básicas para determinar que no faltan datos, es decir, que en las variables originales (salvo el *payload*) no falta ningún tipo de registro. En este caso no hay existencia de *NAs* en las variables, en otras palabras, todos los datos tienen los registros necesarios.

En la siguiente gráfica se puede observar la cantidad de valores únicos en cada variable, en la que se pueden apreciar dos cosas claras:

- Las variables con más únicos son las de los puertos.
- No hay una gran cantidad de IPs únicas.

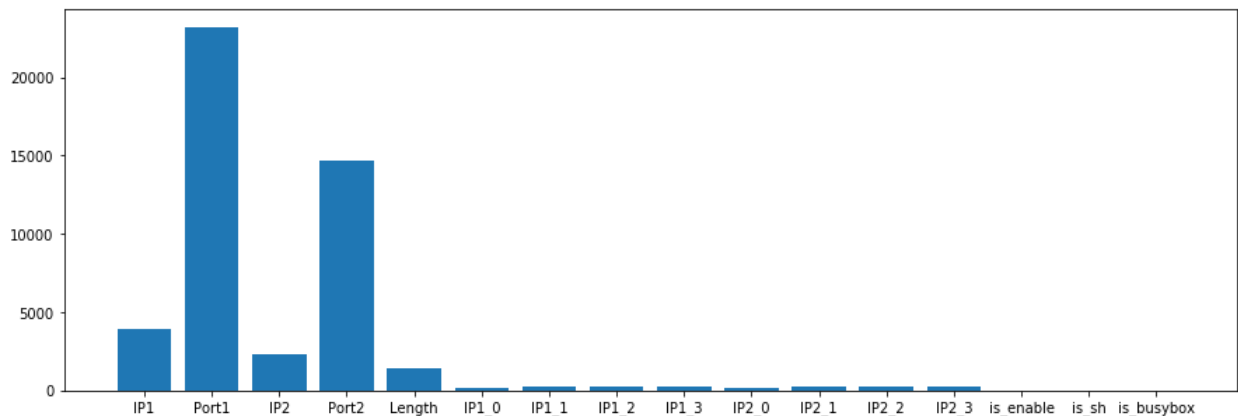


Figura 3.4: Cantidad de valores únicos en las variables

La parte de los puertos en origen (subíndice 1) tiene sentido que sea mayor que en la IP destino, dado que a la hora de iniciar conexiones por lo general el puerto origen es aleatorio y el de destino suele ser fijo, como el puerto 22 para conexiones *ssh*, por lo que puede significar que se estuviera haciendo un escaneo de puertos. Esto se puede ver en el diagrama de barras representando los puertos más utilizados en origen y destino:

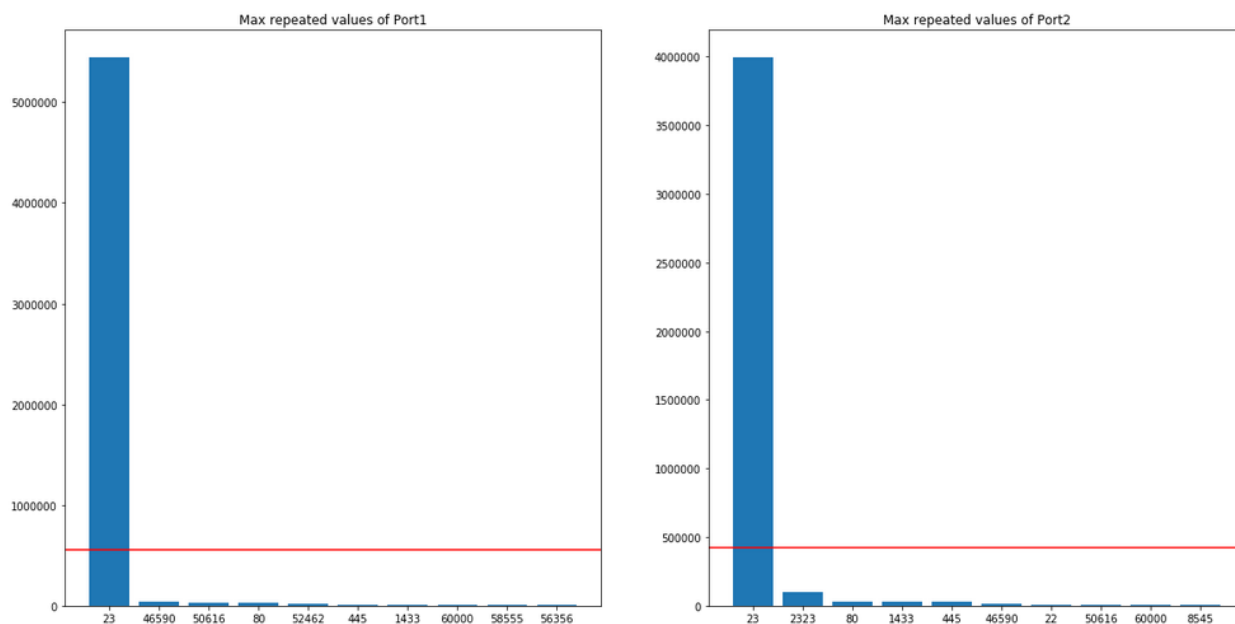


Figura 3.5: Puertos más utilizados

Se puede observar lo siguiente:

- EL puerto más utilizado en ambos casos con diferencia es el puerto 23. Es el puerto utilizado para realizar las conexiones vía *telnet*.
- Seis de los diez puertos más utilizados en la IP de origen son puertos aleatorios, utilizados al realizar la conexión. Los puertos del rango 4xxxx-6xxxx caen en este uso efímero.
- El puerto 80 por lo general es utilizado en conexiones *HTTP* (web).
- El puerto 445 es utilizado para la transferencia de ficheros.
- El puerto 1433 suele utilizarse en conexiones de *Microsoft SQL Server*.
- El 2323 puede ser un puerto alternativo para la conexión vía *telnet*[\[31\]](#).

- El puerto 22 se utiliza para conexiones *ssh*.
- Tanto para destino como para origen se repiten los puertos 46590, 50616 y 60000, cayendo en el rango de puertos temporales.

En el caso de las IPs más utilizadas se puede ver lo siguiente:

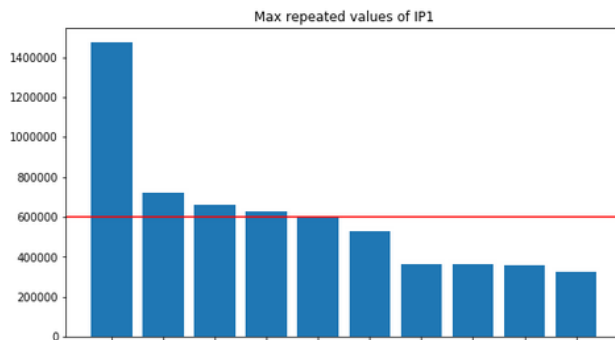


Figura 3.6: IPs de origen con la mayor cantidad de conexiones

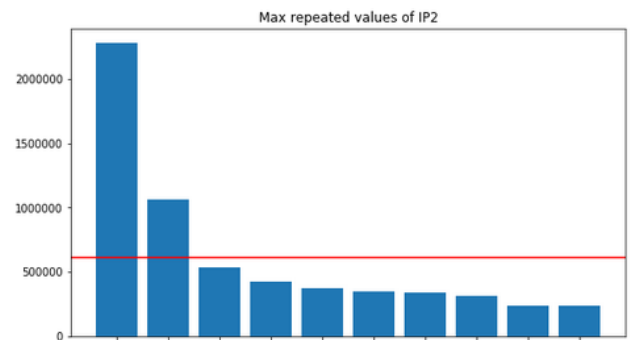


Figura 3.7: IPs de destino con la mayor cantidad de conexiones

Con una cantidad de datos cercana a los 11 millones de registros la mayor parte de las conexiones se registran en estos diez valores máximos con una diferencia de casi el doble de la primera a la segunda tanto para origen como destino.

Dentro de la temporalidad de las conexiones se pueden observar patrones de las conexiones a lo largo del día y su evolución. Debido a la granularidad de los tiempos (capturan con un detalle de microsegundos) estos patrones se han agrupado en distintas ventanas, de modo que en ventanas más pequeñas se puedan detectar ascensos y descensos bruscos, mientras que en las ventanas grandes se puede observar el patrón general. Por ello se han realizado gráficas con ventanas de 1 minuto, 5 minutos y 30 minutos.

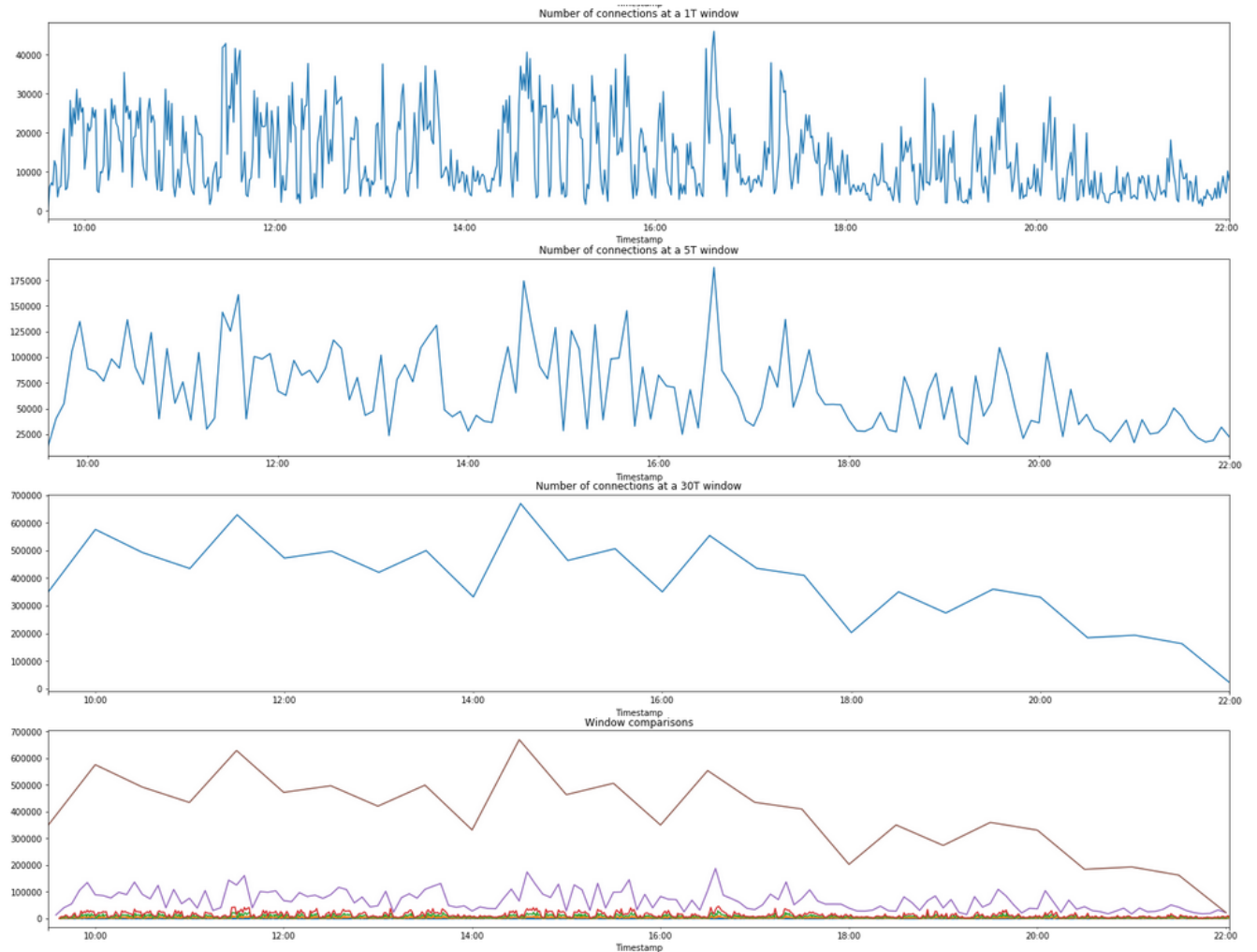


Figura 3.8: Evolución de las conexiones a lo largo del tiempo con distintas ventanas

Se observa:

- Un patrón claro de la cantidad de conexiones, comienza a aumentar durante la mañana y la cantidad de peticiones disminuye por la tarde
- En ventanas más pequeñas existe una caída cerca de las 11 horas y un pequeño pico después.
- Se aprecia un pequeño valle a las 14:00, más pronunciado en ventanas pequeñas.
- Gran caída sobre las 18:00.

- En las comparaciones entre ventanas, se aprecia como el patrón es semejante entre ellas.

En cuanto a las variables binarias, se observa que se encuentran muy descompensadas y existen pocas observaciones que las contengan:

	is_busybox	is_enable	is_sh
0	9812548	10114993	9643741
1	358980	56535	527787

Cuadro 3.1: Tabla de frecuencias de las variables binarias

Por último, en la exploración inicial se observó la distribución de la variable *Length*. En la figura 3.9 se puede apreciar como la mayor parte de las conexiones se realizan con una longitud muy pequeña, recayendo la mayor parte de los valores en torno a 40-50 o con una magnitud grande, generalmente en torno a los 1400, mientras que los valores intermedios son minoría.

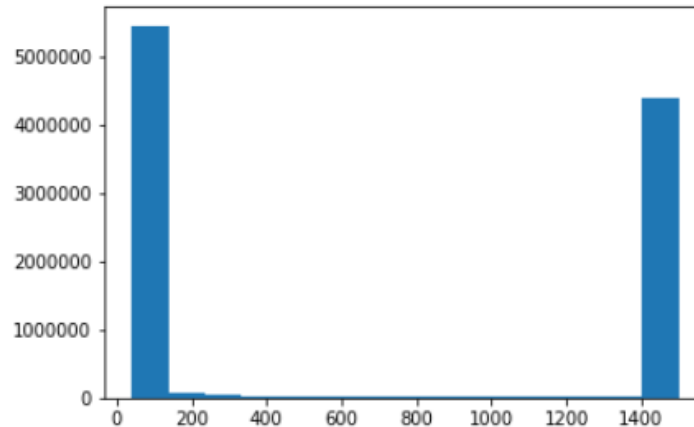


Figura 3.9: Histograma de la variable *Length*

Sobre la variable *Payload* no se realizó ningún tipo de exploración, a parte del parseo, debido a la dificultad de procesar los datos y adaptarlos a un formato más tabular. Se trató de obtener la secuencia de comandos en aquellas observaciones no cifradas sin embargo, no se consiguió obtener esta secuencia que podría haber ayudado en gran medida aportando una gran cantidad de información. Algunas de las variables que se podrían haber obtenido son: cantidad de comandos en secuencia, longitud media de los comandos o incluso técnicas de NLP como *Tf-idf* (*Term frequency – Inverse document frequency*).

Con la información obtenida y con un mayor conocimiento de los datos disponibles, se propusieron nuevas variables que podrían aportar más información y podrían ser más útiles a los modelos (en este momento aún por decidir). Tras varios intentos se generaron las siguientes variables:

- *IP1_time_diff_per_value*: tiempo en segundos, de la diferencia de tiempos desde la anterior conexión y la conexión actual de la IP de origen.
- *IP2_time_diff_per_value*: tiempo en segundos, de la diferencia de tiempos desde la anterior conexión y la conexión actual de la IP de destino.
- *IP1_IP2_time_diff_per_value*: tiempo en segundos, de la diferencia de tiempos desde la anterior conexión y la conexión actual entre IPs únicas de origen y destino.
- *Port1_last_1*: cantidad de puertos únicos utilizados por la IP de origen en una ventana de un minuto
- *Port2_last_1*: cantidad de puertos únicos utilizados por la IP de destino en una ventana de un minuto

Todas estas variables se generan en específico para los valores únicos de las IPs, es decir, cuando se utiliza la IP de origen, se itera entre todos los valores únicos y por cada valor se genera la ventana o la diferencia de tiempo. Con esto se consigue que las nuevas variables solo estén relacionadas con las mismas IPs, tanto de origen, destino o para el caso de *IP1_IP2_time_diff_per_value* de ambas y que no dependan de conexiones entre distintas IPs.

Con las nuevas variables se pretende conseguir distintos tipos de información, con la diferencia de tiempo en la IP de origen, se pretende observar si se está utilizando asiduamente, esta teniendo un pico de conexiones o si por el contrario las establece raramente. Por otro lado, para la IP de destino se pretende observar si esta siendo objetivo de muchas conexiones en cortos periodos de tiempo, que podrían significar ataques *DDoS* (*Distributed Denial of Service*). Para la diferencia entre las conexiones únicas de IP se pretende obtener si se pueden estar realizando escaneos de puertos sobre la IP de destino. En el caso de las variables de puertos, su función es obtener información más precisa de si puede estar realizando un escaneo o existen IPs de origen con muchas conexiones simultáneas que observan distintos destinos.

Una vez generadas estas variables, se necesita observar la cercanía (o distancia) entre las redes, puesto que es muy común que una vez se infecte un dispositivo, éste intente infectar

aquellos dispositivos que se encuentren en su misma red o una red cercana. Inicialmente se planteo utilizar la IP como variable numérica y obtener la distancia euclídea entre la IP origen y destino, sin embargo, no existe ningún tipo de relación entre ellas de esta manera, es decir, la IP 1 y 2 no están a la misma distancia que las IPs 230 y 231. Aprovechando las variables generadas de las IPs con las subredes y el host, se podría calcular la distancia entre ellas, sin embargo, en este caso tomaría como equivalente la distancia entre hosts y la distancia entre redes principales sin embargo, redes con distinto valor no están relacionadas mientras que si la única diferencia es en el host, estos se encuentran en la misma red.

$$dist(IP_{red1}, IP_{red2}) \neq dist(IP_{host1}, IP_{host2}) \quad (3.2)$$

Para ello se tomó la IP como un espacio 4-dimensional (red1, red2, red3 y host) con valores ponderados, 1000 para la primera red, 100 para la segunda, 10 para la tercera y uno para el host.

$$IP_x = (x_1, x_2, x_3, x_4) \quad (3.3)$$

$$dist(IP_x, IP_y) = \sqrt{1000 * (x_1 - y_1)^2 + 100 * (x_2 - y_2)^2 + 10 * (x_3 - y_3)^2 + (x_4 - y_4)^2} \quad (3.4)$$

Con esto una diferencia de uno en la red principal, daría una distancia de $\sqrt{1000} \approx 32$.

Tras realizar pruebas con este tipo de distancia, se decidió utilizar otro enfoque mejorando los resultados. En lugar de tomar la IP como un valor 4-dimensional, se tomo como 32-dimensional, es decir, una dimensión por cada bit.

$$IP = 192,128,1,1 = 110000001000000000000000100000001 \quad (3.5)$$

Para dar aún más peso a los valores iniciales y la cercanía a nivel de red se utilizo el operador *XOR*, devolviendo un 1 para cuando existe diferencia entre los bits de la misma posición o 0 si el valor es igual. Una vez aplicado el operador, se cuentan el número de ceros por la izquierda, en otras palabras, el número de bits iniciales (redes principales) que coinciden, al momento que aparece una diferencia se para de contar. El número resultante se le resta a 32 y se divide entre

32 para normalizar el valor entre 1 y 0, siendo el uno redes completamente distintas y el cero la misma red y el mismo host. Podemos ver el siguiente ejemplo entre IPs que solo se diferencia en el host:

$$\begin{aligned}
 IP_1 &= 192,128,1,1 = 110000001000000000000000100000001 \\
 IP_2 &= 192,128,1,2 = 110000001000000000000000100000010 \\
 XOR(IP_1, IP_2) &= 000000000000000000000000000000010 \\
 n_{zeros} &= 31 \\
 dist(IP_1, IP_2) &= \frac{(32 - n_{zeros})}{32} = \frac{1}{32}
 \end{aligned}$$

En el caso de diferencias entre redes:

$$\begin{aligned}
 IP_3 &= 255,128,1,1 = 111111111000000000000000100000010, \\
 XOR(IP_1, IP_3) &= 00111111100000000000000000000000, \\
 n_{zeros} &= 2, \\
 dist(IP_1, IP_3) &= \frac{(32 - n_{zeros})}{32} = \frac{30}{32}
 \end{aligned}$$

Como se puede observar la distancia aumenta en gran medida con las redes y en muy poca con los host, aportando más valor como nueva variable que las distancias utilizadas inicialmente.

Una vez creadas las nuevas variables, se procede a la exploración de las mismas para su explotación en la creación del dataset final.

Bibliografía

- [1] Rebecca Bace and Peter Mell. Intrusion detection systems. *National Institute of Standards and Technology (NIST)*, 2001.
- [2] Giovanni Vigna and Christopher Kruegel. Host-based intrusion detection. 2005.
- [3] Stefan Axelsson. Research in intrusion-detection systems: a survey. *Department of Computer Engineering, Chalmers University of Technology*, 1998.
- [4] Pedro Garcia-Teodoro, Jesus Diaz-Verdejo, Gabriel Maciá-Fernández, and Enrique Vázquez. Anomaly-based network intrusion detection: Techniques, systems and challenges. *computers & security*, 28(1-2):18–28, 2009.
- [5] Arthur L Samuel. Some studies in machine learning using the game of checkers. *IBM Journal of research and development*, 44(1.2):206–226, 2000.
- [6] Stuart J Russell and Peter Norvig. *Artificial intelligence: a modern approach*. Malaysia; Pearson Education Limited,, 2016.
- [7] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [8] Simon Haykin. *Neural networks*, volume 2. Prentice hall New York, 1994.
- [9] Qiong Liu and Ying Wu. *Supervised Learning*. Springer US, Boston, MA, 2012.
- [10] Gareth James, Daniela Witten, Trevor Hastie, and Robert Tibshirani. *An Introduction to Statistical Learning – with Applications in R*, volume 103 of *Springer Texts in Statistics*. Springer, 2013.
- [11] Frank Rosenblatt. The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, 65(6):386, 1958.

- [12] Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. 96:148–156, 1996.
- [13] Ben Gorman. A kaggle master explains gradient boosting, 2017. Last accessed 27 March 2019.
- [14] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. pages 785–794, 2016.
- [15] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. pages 3146–3154, 2017.
- [16] Home credit default risk, 2018. Last accessed 27 March 2019.
- [17] Corporación favorita grocery sales forecasting, 2018. Last accessed 27 March 2019.
- [18] Google analytics customer revenue prediction, 2019. Last accessed 27 March 2019.
- [19] How to calculate principal component analysis (pca) from scratch in python, 2018. Last accessed 04 April 2019.
- [20] Covariance matrix, 2019. Last accessed 04 April 2019.
- [21] Shyam Boriah, Varun Chandola, and Vipin Kumar. Similarity measures for categorical data: A comparative evaluation. In *Proceedings of the 2008 SIAM international conference on data mining*, pages 243–254. SIAM, 2008.
- [22] Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Kdd*, 1996.
- [23] Anonymous. Inferring and characterizing internet-scale iot probing campaigns by leveraging a novel data dimensionality reduction technique. *ACSAC*, 2018.
- [24] Andrew Ng et al. Sparse autoencoder. *CS294A Lecture notes*, 72(2011):1–19, 2011.
- [25] Analizador de paquetes, 2019. Last accessed 03 June 2019.
- [26] Jon Postel. Transmission control protocol. Technical report, 1981.
- [27] Jon Postel. Internet protocol. Technical report, 1981.
- [28] Feature engineering: What powers machine learning, 2019. Last accessed 04 June 2019.

- [29] Tcpdump/libcap public repository, 2019. Last accessed 04 June 2019.
- [30] Anatomy of an iot malware attack, 2017. Last accessed 04 June 2019.
- [31] What is happening on 2323/tcp?, 2016. Last accessed 04 June 2019.