



UNIVERSIDAD DE MURCIA

TRABAJO FINAL DE GRADO

---

# Identificación de dispositivos IoT a través de huellas hardware

---

*Autor:*

Sergio MARÍN SÁNCHEZ  
sergio.marins@um.es

*Tutores:*

Gregorio MARTÍNEZ PÉREZ  
gregorio@um.es  
Pedro Miguel SÁNCHEZ  
SÁNCHEZ  
pedromiguel.sanchez@um.es

2 de junio de 2022

Me gustaría expresar mi agradecimiento a todos los amigos, profesores y familia que han hecho posible que llegue a este punto.

---

# Índice general

---

<b>Resumen</b>	<b>4</b>
<b>Extended abstract</b>	<b>6</b>
<b>1 Introducción</b>	<b>10</b>
<b>2 Estado del arte</b>	<b>12</b>
<b>3 Análisis de objetivos y metodología</b>	<b>16</b>
3.1 Objetivos . . . . .	16
3.2 Algoritmos de Machine Learning . . . . .	16
3.2.1 Tratamiento de los datos . . . . .	17
3.2.2 Algoritmos de aprendizaje supervisado . . . . .	17
<b>4 Diseño y resolución</b>	<b>21</b>
4.1 Elección del protocolo . . . . .	21
4.2 Obtención de datos . . . . .	21
4.3 Análisis de los datos . . . . .	23
4.3.1 Experimento 1: Muestra secuencial . . . . .	23
4.3.2 Experimento 2: Muestra paralela . . . . .	24
4.4 Elección de la muestra de datos . . . . .	25
4.5 Reducción de la dimensionalidad . . . . .	26
4.6 Entrenamiento de los modelos . . . . .	26
<b>5 Conclusiones y vías futuras</b>	<b>32</b>

---

## Índice de figuras

---

2.1	Desbordamiento del contador [oser2018identifying]	13
2.2	Modelo del sistema de identificación [hamad2019iot]	14
2.3	Clasificador de dos niveles [aksoy2019automated]	15
3.1	Funcionamiento del algoritmo KNN [knnalgorithm2010]	19
3.2	Separación por hiperplanos SVM [svmseparation2012]	20
4.1	Topología de la red	22
4.2	Offset acumulado muestra secuencial	23
4.3	Diferencias entre offsets de dispositivos	23
4.4	Diagrama de cajas muestra secuencial	24
4.5	Offset acumulado muestra paralela	25
4.6	Diagrama de cajas muestra paralela	26
4.7	Correlación entre las variables estadísticas	27
4.8	Particiones de los datos	28
4.9	Comparativa hiperparámetros Random Forest	29
4.10	Comparativa de resultados entre modelos	29
4.11	Matrices de confusión con datos de la muestra paralela	30
4.12	Matriz de confusión del modelo final	31

---

## Índice de tablas

---

2.1	Resultados en el estado del arte . . . . .	15
4.1	Ejemplo de los datos obtenidos de cada dispositivo . . . . .	22
4.2	Datos estadísticos muestra paralela . . . . .	26
4.3	Equivalencia entre algoritmo e implementación . . . . .	27

---

## Resumen

---

En este proyecto hemos diseñado un modelo capaz de identificar dispositivos conectados en red local en base a las diferencias entre sus relojes internos comparados con un reloj que se fija como exacto. Con esto hemos creado un mecanismo que permitirá realizar conexiones más seguras.

Este trabajo se puede dividir en tres partes. Por un lado tenemos la obtención de los datos, para lo que usamos 5 dispositivos iguales en hardware y software (Raspberry Pi). De estos dispositivos obtenemos marcas de tiempo cada segundo de dos formas distintas. Por otro lado, tenemos el análisis estadístico de los datos obtenidos, que generará nuevos datos a partir de los previos. Por último tenemos el desarrollo de un modelo de machine learning que sea capaz de automatizar todo el proceso de distinguir los dispositivos partiendo de los datos estadísticos que hemos obtenido previamente.

Para obtener los datos hemos tenido que realizar una arquitectura cliente-servidor entre los dispositivos a analizar y el dispositivo observador, con el fin de controlar el ritmo de envío de los paquetes. Las capturas se han realizado de dos formas, en secuencial y en paralelo. Por secuencial nos referimos a que primero obtenemos todas las marcas de tiempo de un dispositivo y posteriormente pasamos al siguiente. Por otro lado, en paralelo nos referimos a que se obtienen marcas de tiempo de todos los dispositivos simultáneamente.

Una vez tenemos todas las marcas de tiempo, debemos obtener sus desviaciones respecto al dispositivo que tomamos de referencia. De estas desviaciones se obtiene su incremento entre cada dos puntos, con estos valores usamos una ventana deslizante de 1 minuto, que equivale a 60 paquetes (se toman muestras cada segundo), y con de ellos obtenemos variables estadísticas.

Una vez tenemos las variables estadísticas procedemos a analizar los modelos de machine learning que vamos a usar. Los modelos que se van a analizar son Random Forest, MLP, KNN, Naive Bayes, Árboles de decisión y SVM con kernel lineal.

Para comparar los modelos entre sí usaremos el modelo entrenamiento/validación/test como forma de dividir nuestros datos. Entrenamos todos nuestros modelos con el conjunto de entrenamiento y comprobamos su capacidad de generalización con el conjunto de validación. A la hora de entrenar los modelos comprobaremos diferentes parrillas de hiperparámetros para ajustar nuestros modelos lo más posible a los datos de entrenamiento.

Al realizar estos entrenamientos vemos que los modelos basados en árboles son los mejores, tanto árboles de decisión (97.57 %) como Random Forest (99.51 %). Nos quedamos finalmente

con el modelo de Random Forest como algoritmo final y lo entrenaremos con la totalidad de la base de datos de entrenamiento.

Finalmente el modelo de Random Forest con los hiperparámetros ajustados obtiene un valor de accuracy de 99.44 %.