

# Detección de Phishing mediante Aprendizaje de Máquina

James Garzón Otálvaro  
Estudiante de Ingeniería de Sistemas  
Universidad de Antioquia  
james.garzon@udea.edu.co

Yoiner Esteban Gómez Ayala  
Estudiante de Ingeniería de Sistemas  
Universidad de Antioquia  
yoiner.gomez@udea.edu.co

**Resumen-** El *phishing* es una forma de ingeniería social en el cual el atacante intenta obtener información sensible de una víctima de manera fraudulenta, suplantando una entidad o persona de confianza para realizar acciones o usar recursos del usuario atacado. Siendo el *phishing* uno de los ataques informáticos más frecuentes, es el tema de interés para este artículo donde se usará técnicas de *Machine Learning* (ML) para la predicción de este tipo de ataque en el cual se describe los datos analizados, el contexto en base al estado del arte, los procesos ejecutados y finalmente los resultados obtenidos con sus respectivas conclusiones.

**Palabras Clave-** Phishing Websites, aprendizaje de máquina, seguridad de la información.

## I. INTRODUCCIÓN

El crecimiento de las suplantaciones han aumentado considerablemente en los últimos años [1] y ha retado a los expertos en seguridad informática para ellos encuentren mecanismos que puedan evitar el phishing y de este modo lograr la protección de los datos personales e información sensible de una entidad o compañía. En este caso, se busca predecir los sitios web que tengan como objetivo suplantar la identidad de las personas a partir de un conjunto de datos que provienen de un problema de clasificación biclase por medio de métodos de ML como Funciones Discriminantes Gaussianas, K vecinos más cercanos (K nearest neighbors, KNN), Redes Neuronales Artificiales (RNA), Random Forest y Máquinas de Soporte Vectorial (Support Vector Machines, SVM).

## II. CONJUNTO DE DATOS

### A. Descripción general de la base de datos

El conjunto de datos utilizado para elaborar los experimentos fue tomado de “*Phishing Websites Data Set*” [2] donde se recolectaron 11055 muestras por medio de fuentes como: archivos de *PhishTank*, archivos de *MillerSmiles* y operadores de búsquedas de *Google*. Sus contribuidores han demostrado que el conjunto de datos para realizar las predicciones ha sido satisfactorio respecto a las evidencias reales ya que las 30 características fueron elegidas a partir de unas condiciones comunes que se presentan al momento de efectuar un ataque de este tipo. A continuación se describen las reglas que se utilizaron y el significado de los valores de la salida de cada una de las características.

### B. Descripción de las características

Es importante destacar que las variables se partitionaron en 4 grupos de acuerdo a su relación por parte de los autores de la base de datos.

Tabla I  
CODIFICACIÓN DE CARACTERÍSTICAS

Grupo	Regla condicional	Codificación
$P_1$	Dirección IP o codificación hexadecimal en la URL	0.2
$P_2$	0.3	0.5
$P_3$	0.1	0.2
$P_4$	0.3	0.5

La descripción de las figuras deberá ubicarse debajo de las mismas, centrada, numerándose con cifras arábigas. Use la abreviatura Fig. n tanto para etiquetar la figura o gráfico como para referirse a ella.

La descripción de las tablas deberá ubicarse encima de las mismas, numerándose con cifras romanas y con el texto en versalitas. La etiqueta de la tabla (Tabla X) debe escribirse en mayúsculas y encontrarse sola en una línea. Use Tabla X para referirse a una tabla.

Los pies de las figuras y de las tablas deben seguir el formato mostrado bajo la Fig. 1 y bajo la tabla 1. Si es posible, utilice un formato vectorial (como EPS o PDF) para representar diagramas. Los formatos de tipo *raster* (como PNG o JPG) suelen generar ficheros muy grandes y pueden perder calidad al ampliarlos.

Tabla I  
TABLA DE EJEMPLO

Protocolo 1	Escenario 1	Escenario 2	Escenario 3
$P_1$	0.1	0.3	0.2
$P_2$	0.2	0.3	0.5
$P_3$	0.2	0.1	0.2
$P_4$	0.3	0.3	0.5

### C. Ecuaciones

Las ecuaciones deben estar centradas y situadas en líneas distintas. Cada ecuación debe ser numerada:

$$E = mc^2 \quad (1)$$

Para referenciar una ecuación, utilice Ec. 1.

### D. Numeración, pies y encabezados de páginas

No aplique ningún elemento de numeración, pie o encabezado de página. Estos elementos se añadirán en el proceso final de confección de las actas. Por favor, deje la numeración tal como está en el documento modelo.

### E. Referencias

Las referencias serán numeradas en orden de aparición [1]. El formato de referencias será el estándar del IEEE. Se muestra algún ejemplo en el apartado correspondiente.

### F. Nombre y filiación de los autores

Según el número de autores, adapte la zona correspondiente al nombre y filiación de manera oportuna. Intente no variar de manera notable el aspecto y tamaño de la zona.

## III. CONCLUSIONES

El seguimiento de las normas indicadas permitirá que su trabajo resulte visualmente atractivo. Esta misma plantilla se puede encontrar en formato LATEX, en la dirección *web* oficial de las jornadas (<http://www.jitel.org>).

## REFERENCIAS

- [1] J. Díaz-Verdejo, "Ejemplo de bibliografía", En Actas de las XI Jornadas de Ingeniería Telemática, vol. 1, n. 1, pp. 1-5, 2013.