

Resumen

Con la llegada de Covid-19, han aumentado los correos electrónicos de Spam que intentan engañar a las personas para hacerlas víctimas de estafas o algún engaño. Es posible invertir en software y hardware, pero al final basta con el descuido de un usuario o su desconocimiento para ser víctima de algún ciberataque. Así, este estudio ofrece principalmente dos cosas: primero, hace una revisión de la literatura sobre las soluciones que utiliza Natural Language Processing; y segundo, analiza los resultados obtenidos por los algoritmos Hidden Markov Model (HMM) versus los obtenidos por Machine Learning (ML) en la detección de estos ataques. Con la revisión literaria se puede afirmar que no se han encontrado muchos artículos que utilicen HMM para solucionar este tipo de ataques con una gran precisión; esto se debe a que este tipo de modelo no cuenta con una amplia investigación previa para poder ser aplicado con gran efectividad, sin embargo, su presencia ha abierto una nueva línea de investigación para poder prevenir este tipo de ataques. Para esto, se realizó un estudio comparativo sobre la efectividad y precisión que tienen hoy en día los mejores algoritmos de Machine Learning en cuanto a la detección de ataques por correo y se evidenció que su porcentaje de precisión podría variar según su escenario y la cantidad de datos a procesar. Finalmente, se determinó que los algoritmos de Machine Learning ofrecen mayor precisión en este tipo de detección. Como trabajo futuro se propone el desarrollo de un algoritmo que realice un preprocesamiento mediante HMM y luego utilice algoritmos de Deep Learning para mejorar la precisión.

Palabras clave: Spam, correo electrónico, Hidden Markov Model, Machine Learning, Ingeniería Social.

Abstract

With the arrival of Covid-19, there has been an increase in Spam e-mails that try to trick people into becoming victims of scams or deception. It is possible to invest in software and hardware, but in the end a user's carelessness or lack of knowledge is enough to become a victim of a cyberattack. Thus, this study offers mainly two things: first, it makes a literature review on the solutions that use Natural Language Processing; and second, it analyzes the results obtained by Hidden Markov Model (HMM) algorithms versus those obtained by Machine Learning (ML) in the detection of these attacks. With the literature review it can be stated that not many articles have been found that use HMM to solve this type of attacks with high accuracy; this is because this type of model does not have extensive previous research to be applied with great effectiveness, however, its presence has opened a new line of research to prevent this type of attacks. For this, a comparative study was carried out on the effectiveness and accuracy that the best Machine Learning algorithms have nowadays regarding the detection of mail attacks and it was evidenced that their percentage of accuracy could vary according to their scenario and the amount of data to be processed. Finally, it was determined that Machine Learning algorithms offer higher accuracy in this type of detection. As future work, we propose the development of an algorithm that performs preprocessing using HMM and then uses Deep Learning algorithms to improve accuracy.

Key words: Spam, email, Hidden Markov Model, HMM, Machine Learning, Social Engineering.