# AI in Cyber-security: Spam Filtering

JAHJA Darwin, 16094501D

Artificial intelligence (AI) has become one of the biggest trends in recent years. This new technology has been widely adopted in many different application domains such as logistics and marketings, and one of its most important integrations is towards cyber-security.

Cyber crimes happen everyday in our real life through e-mail, phone call, social media etc., and this trend keeps evolving and may threaten people's lives as well as their properties. To fight against cyber criminals, spam filter has been invented to help prevent malicious and unsolicited information reaching to potential victims. With the integrations of AI, modern spam filter has evolved and become more effective and powerful in detecting and filtering out suspicious and junk messages before they reach to the victims.

This report will mainly discuss the AI's contributions to cyber-security by studying different practical applications of AI-integrated spam filters.

#### E-mails

E-mail services have been commonly used not only for personal communication, but also for commercial advertisement in today's world. The high usage of e-mail has stimulated several cyber-security threats caused by spam e-mails, which include viruses, phishing messages, fraud schemes, explicit contents etc. (Bhowmick & Hazarika, 2018). To address these issues, companies and E-mail Service Providers (ESP) create and customize their own e-mail spam filters to block spam e-mail. Google Mail (GMail), one of the largest ESPs, has integrated AI technology into its spam filter to protect its users from being victimized by cyber criminals (Metz, 2018).

The techniques used by e-mail spam filter can be categorized into two different types: the traditional SMTP approach and the newer machine learning approach, which the later uses machine learning classifier to detect whether the e-mail is spam or legitimate. It is consider to be more efficient and has more popularity as it can learn the patterns and the changes of the new spamming techniques automatically time by time, without manually configuring by humans (Roy & Viswanatham, 2016).

The machine learning spam filtering techniques can be further subcategories into content based and non-content based. For content based, the machine learning algorithm used varies base on its type. For example, Naive Bayes Classifier or Support Vector Machine (SVM) is used in text based contents to filter contents such as phishing messages and fraud schemes (Roy & Viswanatham, 2016), while Back Propagation Neural Networks is used in image base contents to filter explicit contents such as sexual images (Chowdhury, Gao & Chowdhury, 2015). For non-content based, N-Grams Algorithm is used to classify the e-mail based on the e-mail header, which consist of information such as the sender server IP address and the mail subject (Hu et al., 2010). This can prevent e-mail abusing behaviors such as spambot attack and e-mail bombing attack.

With the integration of AI into e-mail spam filter, not only it can protects its users against cyber criminal, but it can also learn from the user's preferences and customize the received content. The aim is to provide a better spam-free, user-friendly e-mail communication experience for people.

#### Voice Calls

Similar to e-mail, voice call is widely used for both personal communication and commercial activity, and it attracts cyber criminals to perform illegal activities such as scamming using the internet's information. As people implicitly or explicitly give out their personal information (PI) on the internet, cyber criminals may have the chance to utilize their PI and perform attacks (Liu et al., 2015). To solve this problem, besides educating the public to protect their PI, spam call filter has been made to block nuisance callers. Nowadays, companies such as Google and DialogTech build AIs that utilize machine learning and artificial neutral networks to detect and screen spam call (Tillman, 2018).

The techniques used to classify spam calls are similar to that in classifying spam e-mails. In addition to that, real time *Natural Language Processing* (NLP) is used to process the conversation in a call, which the *Automated Speech Recognition* and *Natural Language Understanding* in NLP convert the caller's speeches into text based dialogs, and then used by the machine learning classifier to identify whether the call is spam or not. (Li et al., 2018)

Moreover, based on these techniques, spam call filter can also distinguish human and robot speakers. Taking Google Call Screen as an example, when the AI detect a robocall, it will keep dangling so that the robot may not be able to make a new call unless the current call is ended (Orlowski, 2018). This can eliminate the chances for nuisance robocalls wasting people's resources and thus provide a spam-free experience.

### Social Media

Besides e-mail and voice call, the majority of people nowadays use social media to communicate and share their lives and feelings with others through posts and comments. In the past, spamming was one of the biggest issues that all social media face as it is difficult to detect and eliminate spam contents (e.g. fake news, hate speech, nudity, harassment etc.) automatically (Fuchs, 2017). In order to solve this problem, popular social media such as Facebook and Instagram have developed powerful AIs to tackle spam contents (Roettgers, 2017).

The techniques they have used is called Deep Learning, which enables computers to have deeper and better understanding to texts, speeches or images just like human does (Abdulkader et al., 2016). By using data generated from different users, the models can be successfully trained and applied into real life. For example, both Facebook and Instagram use *DeepText*, a text understanding engine developed by Facebook, to build their comment and spam filters. Based on this technology, the spam filter can further understand the meaning of the posts and comments and determine the kinds of content (Abdulkader et al., 2016). Facebook also uses its *Image Recognition AI* to distinguish spam and abuse images. Recently, they are collecting nude photos to develop a model that will filter out nudity contents (Solon, 2017).

# **Future Development**

In my opinion, the future of AI integrated spam filter is bright as computers can automatically learn to detect new kinds of spams without any human efforts. Besides filtering out the spams, another effective way to fight against cyber criminals such as scammers is to tackle them back using AIs. The idea is to waste their time with a never ending conversation with an AI so that they will have less time to pursue real people, and eventually forcing them to give up their illegal behavior. For examples, NetSafe has developed Re:scam, an AI e-mail bot that have many personas, to tackle against email scammers (NetSafe, 2017), while Jolly Roger Telephone Company develops an AI calling bot to fight against evil telemarketers (Bilton, 2016).

In conclusion, AI technologies open potentials to the future of cyber-security. AI-integrated spam filter is one of the greatest examples that demonstrates the effectiveness of AI technologies being used in cyber-security. As criminal's techniques evolve time by time, so do AIs evolve and adapt with the new changes. In this Big Data era, I believe that AIs can perform better in protecting us in the cyber world.

## Reference

Bhowmick A., Hazarika S.M. (2018) E-Mail Spam Filtering: A Review of Techniques and Trends. In: Kalam A., Das S., Sharma K. (eds) Advances in Electronics, Communication and Computing. Lecture Notes in Electrical Engineering, vol 443. Springer, Singapore

Roy, S. S., & Viswanatham, V. M. (2016). Classifying Spam Emails Using Artificial Intelligent Techniques. International Journal of Engineering Research in Africa, 22.

Hu, Y., Guo, C., Ngai, E. W. T., Liu, M., & Chen, S. (2010). A scalable intelligent non-content-based spam-filtering framework. Expert systems with applications, 37(12), 8557-8565.

Chowdhury, M., Gao, J., & Chowdhury, M. (2015, October). Image spam classification using Neural Network. In International Conference on Security and Privacy in Communication Systems (pp. 622-632). Springer, Cham.

Metz, C. (2018). Google Says Its AI Catches 99.9 Percent of Gmail Spam. Retrieved from https://www.wired.com/2015/07/google-says-ai-catches-99-9-percent-gmail-spam/

Liu, Y., Song, H. H., Bermudez, I., Mislove, A., Baldi, M., & Tongaonkar, A. (2015, November). Identifying personal information in internet traffic. In Proceedings of the 2015 ACM on Conference on Online Social Networks (pp. 59-70). ACM.

Tillman, M. (2018). What is Google Call Screen and how does it work?. Retrieved from https://www.pocket-lint.com/apps/news/google/146018-google-call-screen-how-to-screen-spam-calls-with-google-assistant

Orlowski, A. (2018). Google offers to leave robocallers hanging on the telephone. Retrieved from https://www.theregister.co.uk/2018/07/10/google\_offers\_to\_leave\_robocallers\_dang ling\_on\_the\_telephone\_line/

Li, H., Xu, X., Liu, C., Ren, T., Wu, K., Cao, X., ... & Song, D. (2018). A Machine Learning Approach To Prevent Malicious Calls Over Telephony Networks. arXiv preprint arXiv:1804.02566.

Fuchs, C. (2017). Social media: A critical introduction. Sage.

Roettgers J. (2017). Instagram Starts Using Artificial Intelligence to Moderate Comments. Is Facebook Up Next?. Retrieved from https://variety.com/2017/digital/news/instagram-ai-machine-learning-facebook-filters-1202482031/

Abdulkader, A., Lakshmiratan, A., & Zhang, J. (2016). Introducing DeepText: Facebook's text understanding engine. Retrieved from https://code.fb.com/ml-applications/introducing-deeptext-facebook-s-text-understanding-engine/#

Solon, O. (2017). Facebook asks users for nude photos in project to combat 'revenge porn'. Retrieved from https://www.theguardian.com/technology/2017/nov/07/facebook-revenge-porn-nude-photos

NetSafe. (2017). Re:scam. Retrieved from https://www.netsafe.org.nz/rescam/

Bilton, N. (2016). A Robot That Has Fun at Telemarketers' Expense. Retrieved from https://www.nytimes.com/2016/02/25/fashion/a-robot-that-has-fun-at-telemarketers-expense.html