JOURNAL OF INFORMATION PROCESSING SYSTEMS JIPS

# A Deep Learning Approach for Identifying User Interest from Targeted Advertising

Wonkyung Kim[1], Kukheon Lee[1], Sangjin Lee[1], and Doowon Jeong[2,*]

### Abstract

In the Internet of Things (IoT) era, the types of devices used by one user are becoming more diverse and the number of devices is also increasing. However, a forensic investigator is restricted to exploit or collect all the user's devices; there are legal issues (e.g., privacy, jurisdiction) and technical issues (e.g., computing resources, the increase in storage capacity). Therefore, in the digital forensics field, it has been a challenge to acquire information that remains on the devices that could not be collected, by analyzing the seized devices. In this study, we focus on the fact that multiple devices share data through account synchronization of the online platform. We propose a novel way of identifying the user's interest through analyzing the remnants of targeted advertising which is provided based on the visited websites or search terms of logged-in users. We introduce a detailed methodology to pick out the targeted advertising from cache data and infer the user's interest using deep learning. In this process, an improved learning model considering the unique characteristics of advertisement is implemented. The experimental result demonstrates that the proposed method can effectively identify the user interest even though only one device is examined.

### Keywords

Convolutional Neural Network (CNN), Deep Learning, Digital Forensics, User Interest, User Profiling

# 1. Introduction

Digital forensics aims to gain a better understanding of an event of interest by finding and analyzing the facts related to the event [1]. To reveal the truth of an event, the digital forensic investigator should examine the remnants of an event left on the digital system, like desktop, laptop, smartphone, wearable device, etc. [2].

As information and communication technology has developed, a variety of portable devices such as Internet of Things (IoT) devices have been released. There was an attempt to leverage IoT devices in digital investigation [3]. However, the rapid development of IoT raises new challenges in digital forensics [4,5]. It is becoming nearly impossible to investigate entire devices or storages due to technical and legal issues.

The increase in storage capacity has required extensive resources and time, therefore, even a simple forensic operation like the search for a target file can be a considerable expense [6]. There are also jurisdictional legal issues as the digital devices may be out of the site or be dispersed in multiple countries

[7]; most used applications use cloud or fog computing that stores files or data on not an individual device but a server system [8]. Furthermore, seizing and analyzing all devices of the suspect can raise an invasion of privacy [9]. For these reasons, the forensic investigators should find out, as much as they can, the event by analyzing the remnants left in some of the suspect's devices.

In this paper, we focus on targeted advertising stored in user's local devices. When the user utilizes a web browser, online platform services analyze the user's usage pattern to extract the user's interest. Because people usually log in to multiple devices with the same account (e.g., Google, Apple, Microsoft), by analyzing the targeted advertising stored in seized devices, keywords or interests that the user has searched for on other devices can be identified.

The remainder of this paper is organized as follows. In Section 2, the background is introduced. Section 3 describes our proposed methodology and then it is proven by an experiment in Section 4. Concluding remarks are drawn in Section 5.

# 2. Background

## 2.1 Deep Learning

The concept of deep learning and the backpropagation algorithm first appeared in the 1980s. At that time, although the algorithm worked successfully, it took too much time to train a neural network. So deep learning was considered as an impractical technique for general applications in other fields. However, in the mid-2000s, with the development of computer hardware and technology that can overcome the existing problems of deep learning (i.e., vanishing gradient problem, overfitting, and initialization), deep learning has begun to draw attention again. Conventional machine learning techniques were limited in their ability to process natural data in their raw form. Constructing the machine learning system requires careful engineering and considerable domain expertise to design a feature extractor that transforms raw data into a suitable internal representation or feature vector [10]. Unlike machine learning, a deep learning method is a set of methods that allow a machine to be fed with raw data and to automatically discover rules needed for detection or classification [11].

## 2.2 Natural Language Processing

Natural language processing (NLP) is a set of technologies that analyze, extract, and understand meaningful information from human text. There are many areas where NLP is used in our daily life; text summary, conversation system (e.g., Apple Siri, Amazon Alexa), and machine translation (e.g., Google Translate) [12]. Most of the early NLP techniques used a rule-based approach to write down human vocabulary and language. However, due to the variability, ambiguity, and context-dependent interpretation of the human language, the conventional NLP technique had a big limit on its performance, and most of the state of art NLP technique has used machine learning for the sophisticated linguistic analysis.

The state of art NLP technique can identify syntactic and semantic information including contextual information [13]. Word embedding is one kind of NLP technique. The word embedding means a process of converting natural language into a vector. Embedding plays a key role as the first gateway for the computer to understand natural language. When deep learning and NLP are used together, the quality of the embedded vector largely affects the overall performance of the whole model.

## 2.3 Targeted Advertising

Online advertising is an enormous industry. The market share of Internet advertising has been steadily rising, currently totaling a third of the total advertising revenues in the United States, and almost equals that of broadcast television [14]. The targeted advertisement is valuable to advertisers as it increases the probability that the advertisement leads to a customer's purchase. Recent research shows that personalized advertisement is twice as effective as non-personalized advertisement [15]. Companies try to gather, process, and manage the information for its optimum use. Traditionally, tracking web cookie was the main way of gathering customer information. But due to the increased usage of mobile devices, browser fingerprinting is becoming a popular way of user identification in place of the traditional cookie-based approach [16]. Personal information used for targeting includes keywords entered in search engines, recent browsing history, previous web purchases, and even the topics in their e-mails. It goes without saying that targeted advertising can be used to infer the user's interest [17].

## 2.4 Related Works

Estimating user interests is a very important issue in areas such as digital forensics and recommend systems. Recently, with the prevalence of various IoT devices, the frequency of users using Social Network Service (SNS) has increased [18]. Accordingly, most studies related to user interests have been focusing on articles written by users on SNS. Kang and Lee [19] try to extract user interests by applying frequency-based NLP techniques to the articles. However, this method may have problems with accuracy. Since it is based on the user-written text, it cannot extract the user's interest if the user-written text does not reflect the user's interests or personality.

Particularly in the digital forensics field, it is a very important issue to extract user interests. As the capacity of storage devices becomes larger and one user uses multiple IoT devices, the importance of user profiling through user interest extraction becomes a very important issue. Forensic researchers have tried to identify user interest by using data collected from web browsing behavior. Luo et al. [20] proposed a method to estimate the user interest by analyzing the user's action on the web. Diep et al. [21] proposed an unsupervised method to estimate the user interest from browsing behavior data. Siriaraya et al. [22] also estimated user interest by analyzing the browsing history of the user. However, because the previous works focused on data stored on local devices, there was a significant drawback that it is difficult to identify how the user utilizes other devices. To overcome the drawback, we estimate the user interest using web cache data of targeted advertising image, which is generated automatically and shared between devices through account synchronization.

Also, looking at the recently published papers, there are many attempts to apply NLP to deep learning. Karbab et al. [23] used deep learning and NLP to determine whether a file is malicious and to identify malware clusters. They performed word embedding for the behavioral report of malware and classified the embedded vector into two classes (benign, malicious). The proposed framework showed F1-score that exceeded 90% on average. In the same manner, Lee et al. [24] proposed a system that detects fake news using deep learning and NLP. They embedded the head and body of the news respectively using fastText and then classified the vector using deep learning. The proposed system showed an accuracy of about 70% for 100K test datasets. Salminen et al. [25] proposed a framework that determines whether comments on SNS platforms such as YouTube, Reddit, Wikipedia, and Twitter are harmful or not. The proposed framework showed an F1-score of about 90% for 200K datasets.

# 3. Methodology

In the case of a normal image, the subject of the image can be extracted exactly just by classifying the image using deep learning (e.g., convolutional neural network). However, in the case of advertising images, even in advertising images that promote the same product, the texts or objects in the images may be completely different (Fig. 1). Considering this characteristic, the existing approach for image classification is not suitable for advertising images. So, we propose a new method that uses an embedded vector extracted from the advertisement to classify the image instead of classifying the image directly.



|      (a)      |      (b)      |

**Fig. 1.** (a, b) Examples of two targeted advertising images containing the same subject (beverage).

Table 1 represents the result of image classification for the two examples using naive deep learning image classifier. Xception is a kind of deep learning image classification model that was implemented by Google in 2016. Google Cloud Vision is a deep learning-based image processing API that is commercially provided by Google. There are three modules in Google Cloud Vision to handle image: label detection, object detection, and text detection. We confirmed that among the modules, label detection is the most suitable for extracting the subject of the image. Therefore, we compared the results of label detection of Google Cloud Vision and Xception.

**Table 1.** Results of applying naive deep learning image classifiers to the two sample images of Fig. 1

|                     | Fig. 1(a)                          | Fig. 1(b)                                |
|---------------------|------------------------------------|------------------------------------------|
| Xception            | Pop bottle, Beer bottle, Wine bottle | Refrigerator, Eggnog, Pop bottle         |
| Google Cloud Vision | Bottle, Drink, Glass bottle        | Vintage advertisement, Drink, Soft drink |

As seen in Table 1, the label detection showed slightly better results than Xception, but all of the two naive deep learning models could not find out that Fig. 1(a) and 1(b) belong to the same category. For this reason, we used deep learning along with word embedding to extract the subject of advertising images.

Our proposed method consists of four major parts: targeted advertising image extraction, subject of image extraction, word embedding, and classification using deep learning. The overall configuration of the proposed system is represented in Fig. 2. The detailed description of each major part is as follows.

**Phase 1. Extracting targeted advertising**

The web browser's cache contains downloaded images, videos, documents, and executable files, so forensic researchers have studied how to use the cache as evidence [26,27]. Our proposed method focuses on Chromium web browser that stores the cache data in "%UserProfile%\AppData\Local\Google\Chrome\Userdata\Default\Cache." Among the cache data, only the targeted advertising images are extracted based on certain rules such as image URL or image size.

**Phase 2. Extracting subject of image**

The subject extraction module extracts the subject of the targeted advertising image. Generally, it is known that the subject of an image can be derived by using object detection or text detection [28]. However, because the advertisement has different characteristics from the general image, we use label detection provided by Google, which conducts a comprehensive analysis and shows slightly better performance than the other naive deep learning model.

**Phase 3. Embedding image as a vector**

This module embeds the derived subjects of the image into a vector. We use fastText developed by Facebook AI Research Lab for word embedding [29]. Because the number of subjects is not enough to train a new fastText model, this module uses a pre-trained model that 2 million-word vectors were trained by web crawling.

**Phase 4. Training system**

A deep learning model is implemented to classify the calculated vectors according to the image categories. This module uses the Keras library to implement the deep learning model. Keras is an open-source neural network library developed in Python [30]. It is a high-level library of TensorFlow for building and training deep learning models and is now used in many areas due to its user-friendly design.
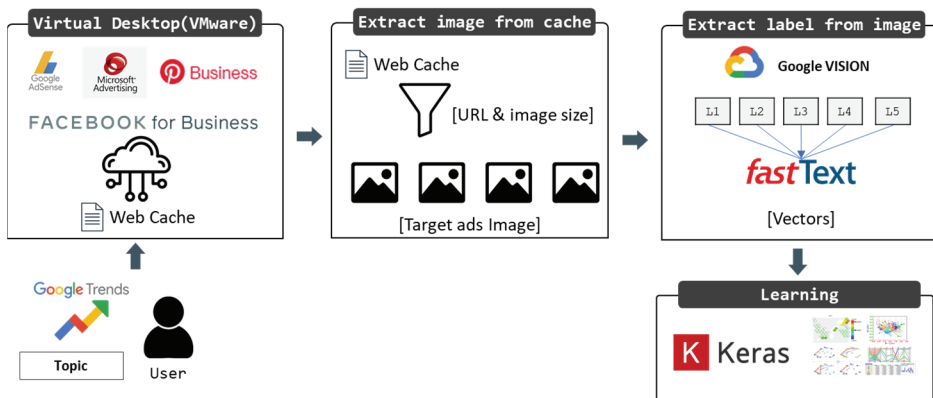


**Fig. 2.** The overall configuration of the proposed system.

# 4. Experiment

## 4.1 Experimental Setup

The experiment was conducted using an i7-8700 processor and an NVidia GeForce 1070 Ti graphic card. First of all, we selected six categories that can represent user interests using Google Trends:

Automobile, Beverage, Clothing, Cosmetic, Electronic, and Food. After selecting the six categories, we chose the detailed keywords on Wikipedia Hierarchy for each category. Then, we repeatedly searched the chosen keywords on websites such as Amazon and Google until the user's interests are reflected in targeted advertising.

Secondly, we extracted targeted advertising images from web cache data based on the size of the image and URL that represents a source of the image. By parsing the cache image, the URL recorded by the advertising agency can be identified. For example, the cache image created by Google AdSense, the most used custom advertising engine, contains the URL that starts with "tpc.googlesyndication." Also, source URL such as "cdn-aitg.wideplanet" and "pix.as.criteo" is recorded in the targeted advertising image. Besides, we have empirically confirmed that the image that size is 4 kB or less is not the advertising image but the icon cache, so we excluded the images below 4 kB. We repeated the above process for cache data, and, as a result, collected experimental data as seen in Table 2.

**Table 2.** Number of collected images by category

| Category | Number of images |
|---|---|
| Automobile | 4,532 |
| Beverage | 4,564 |
| Clothing | 5,323 |
| Cosmetic | 4,905 |
| Electronic | 2,886 |
| Food | 3,076 |

And then, we derived the subjects of the image using label detection of Google Cloud Vision for all the collected images. Among the extracted labels (subjects), the top 5 labels by the score were embedded as 300-dimension vectors using a pre-trained fastText word embedding model. Then, we calculated the weighted average of the five vectors; we consider the weighted average vector as a representative vector for each image.
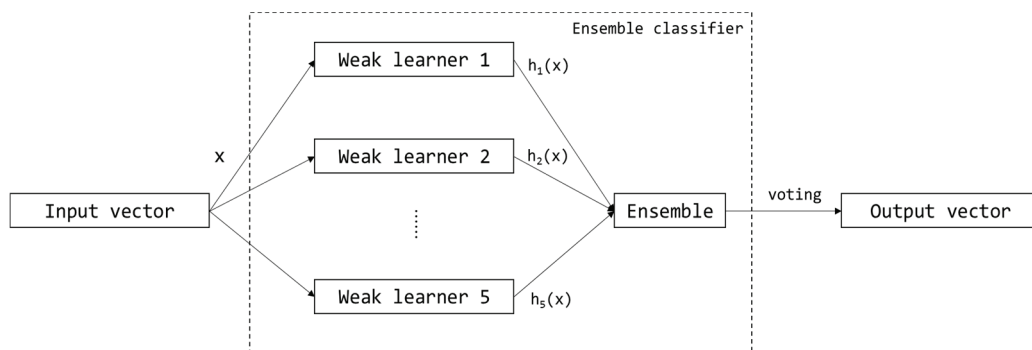


**Fig. 3.** Structure of the ensemble model.

Finally, we used deep learning to classify the representative vector into six categories. We implemented three different deep learning models and compared the performance of each model: the multi-layer perceptron (MLP) model, convolutional neural network (CNN) model, and recurrent neural network

(RNN) model. Additionally, a voting-based ensemble technique using five weak learners was used to improve the performance of each classification model. Ensemble technique is one of the most popular and effective techniques in deep learning. It is a technique that combines and uses several individually trained models that are expressed as weak learners. We built three ensemble models by training five naive models for each of the three naive model types (MLP, CNN, RNN) mentioned above and comparing the results. The use of the ensemble technique improved the overall performance of the model and had the effect of alleviating overfitting to some extent. The structure of the ensemble model we used in the experiment is shown in Fig. 3.

## 4.2 Experimental Result

Our proposed system classified advertising images with an F1-score of 79%. Table 3 shows the overall classification performance and hyperparameters of the model. As seen in Table 3, the model with the ensemble technique shows a performance improvement of about 2% compared to the model without the ensemble. Table 4 is a classification report showing an overall performance of the CNN ensemble model that showed the best performance in our experiment. The confusion matrix of the CNN ensemble model is represented in Table 4.

**Table 3.** Metrics of naive model and ensemble model

| Model | Epochs | Learning rate | Accuracy (%) | |
|---|---|---|---|---|
| | | | Naive | Ensemble |
| MLP | 10 | 0.01 | 74.47 | 76.30 |
| CNN | 10 | 0.01 | 75.98 | 76.69 |
| RNN | 20 | 0.01 | 74.91 | 76.06 |

**Table 4.** Classification report of CNN ensemble model

| | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| Automobile | 0.94 | 0.91 | 0.92 | 452 |
| Beverage | 0.72 | 0.78 | 0.75 | 444 |
| Clothing | 0.74 | 0.84 | 0.79 | 547 |
| Cosmetic | 0.82 | 0.75 | 0.78 | 452 |
| Electronic | 0.68 | 0.70 | 0.69 | 290 |
| Food | 0.88 | 0.71 | 0.79 | 330 |
| Accuracy | - | - | 0.79 | 2515 |
| Macro avg. | 0.80 | 0.78 | 0.79 | 2515 |
| Weighted avg. | 0.80 | 0.79 | 0.79 | 2515 |

In our experiment, the CNN model with the ensemble technique showed the best overall performance. We extracted subjects of the targeted advertising image and then embedded the subjects into 300-dimension vector. Since the meaning of subjects extracted from the image is mostly similar and embedded vector contains the meaning of the word, vectors generated by words with similar meaning have a locally similar distribution. Due to the nature of CNN that specializes in capturing the most salient local features [31], it shows the best performance in our experiment. The model architecture of our CNN model is represented in Fig. 4.
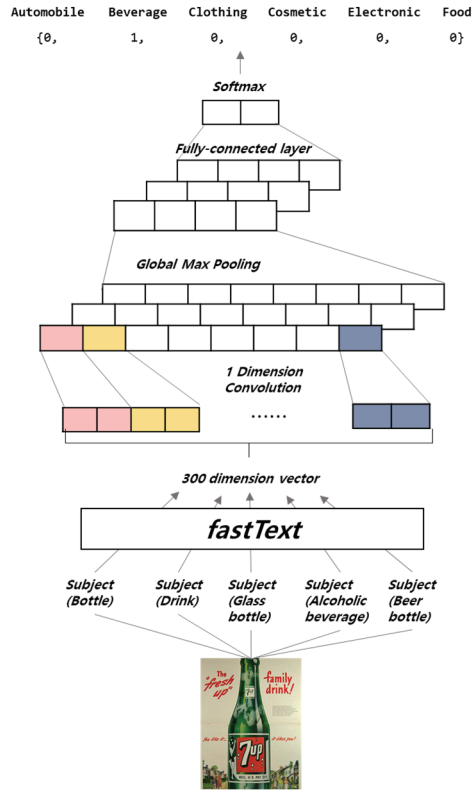
**Fig. 4.** The model architecture of CNN model.
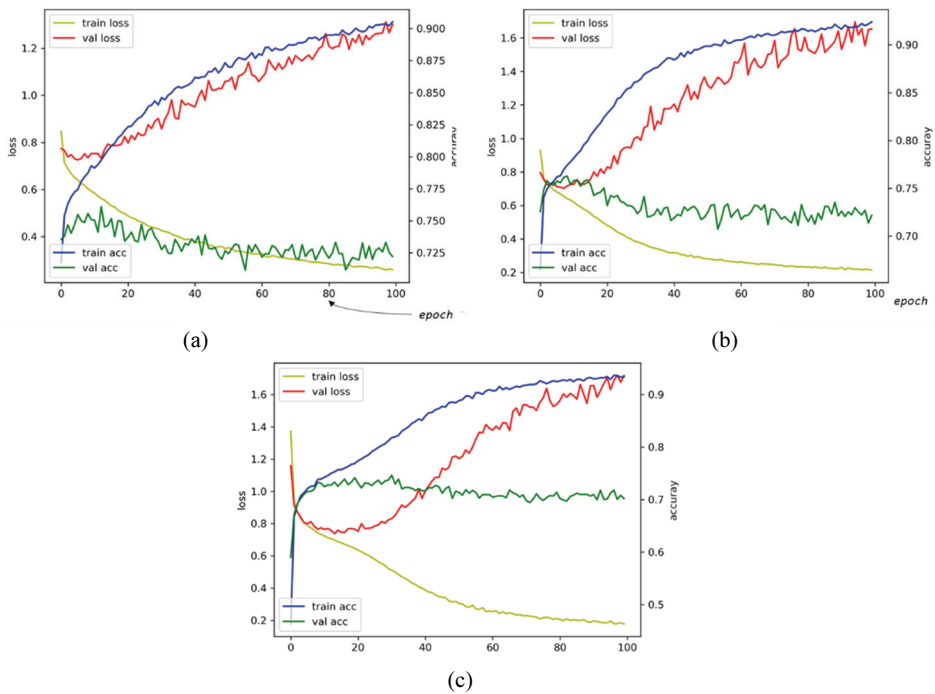


(a)

(b)

(c)

**Fig. 5.** The loss and accuracy of the proposed models: (a) MLP, (b) CNN, and (c) RNN.

The loss and accuracy during training (100 epochs) are shown in Fig. 5. As shown in Fig. 5, MLP and CNN show overfitting when the epoch exceeds 10. In RNN, overfitting occurs when the epoch exceeds 20. In this experiment, because the number of training data is not enough to train a complex model fully, the model tends to occur overfitting quickly.

## 4.3 Discussion

As experiment results have shown, there are some false positive and false negative errors in the system we proposed. We discuss the results as follows:

1) The small dimension of the input vector results in constraints on performance improvements. We could not train a new word embedding model since the number of images collected was too small, so we used a pre-trained word embedding model. As a result, the dimension of the input vector fed into the deep learning model was 300. Because the dimension of the input vector was small, the total number of trainable variables was limited, and it caused insufficient performance.

2) The number of training data is too small and unbalanced. Since the number of training data was too small, we could not build a delicate model. We built a shallow model that combines network layers simply and this results in early overfitting and low accuracy. Also, as seen in Tables 4 and 5, the precision, recall, and F1-score are significantly low in the electronic category. We can deduce the reason for this as an imbalance of the number of training data. In Table 2, it can be seen that the number of learning data in the electronic category is remarkably small. This might cause poor classification performance for the electronic category.

3) It is difficult to extract the subject in case of an advertising image containing content that is not related to the subject of the advertisement. Also, if the result of embedding subjects extracted as a result of label detection does not have a similar distribution in the vector space, this can have a negative influence on the training process.

4) Because we use a pre-trained fastText model, the performance of the proposed system is dependent on the performance of the model. In other words, the quality of the pre-trained fastText word embedding model affects the overall performance of our model. Since the pre-trained fastText model that we used in our experiment was developed not for keywords related to commercial advertisements, but for general purpose. Therefore, there is a limitation in improving the performance of our system.

**Table 5.** Confusion matrix of CNN ensemble model

| Actual class | Predicted class | | | | | |
|---|---|---|---|---|---|---|
| | Automobile | Beverage | Clothing | Cosmetic | Electronic | Food |
| Automobile | 408 | 8 | 13 | 4 | 22 | 5 |
| Beverage | 3 | 356 | 38 | 16 | 18 | 16 |
| Clothing | 14 | 24 | 453 | 16 | 23 | 10 |
| Cosmetic | 3 | 43 | 53 | 370 | 22 | 2 |
| Electronic | 5 | 19 | 27 | 10 | 209 | 3 |
| Food | 2 | 40 | 18 | 4 | 21 | 217 |

Although there are some drawbacks, the overall F1-score of our system is measured as 79%. The experiment confirms that user interest can be identified from advertising images using deep learning and NLP. The targeted advertising image may be found on various devices using the synchronization between

accounts, so forensic investigators can extract the user interest remaining in other devices that cannot be seized, by only analyzing the targeted advertising on the seized devices (Fig. 6).
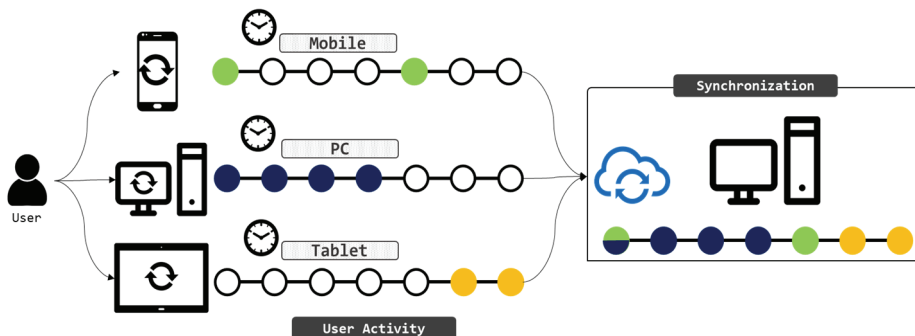


**Fig. 6.** Identification of the information of the user interest scattered in various devices.

# 5. Conclusion

In this paper, we have proposed a system that extracts the user interest from the targeted advertising image in web cache data stored in various devices. As the user and the devices have a close relationship, the user's personal characteristics such as user interests are reflected in the system. As a result, the user's devices become a rich source of useful information for profiling users. The existing methods of extracting user interests using SNS posts cannot accurately extract user interests unless the user reflects the user's characteristics in the posts. And, simple frequency-based analysis on the SNS posts cannot extract the user interests that change over time. Another analysis method of user interests using web browser data, which have been widely used in the digital forensics field, has a disadvantage in that data stored in other devices cannot be analyzed.

The key concept of our proposed system is to leverage web artifacts synchronized between the user's devices. We have used the advanced NLP and deep learning technique to build a classification model for the targeted advertisements which is stored in synchronized web data. The proposed system achieves over 79% F1-score in identifying user interest. Through the experiments, we have confirmed that our proposal is suitable for extracting user's interests in digital forensic practice.

Since the user's interests that have changed over time are reflected in the targeted advertising and updated constantly, the user's interests can be accurately extracted by using this. Also, considering that people usually log in to multiple devices with the same account (e.g., Google, Apple, Microsoft), keywords or interests that the user has searched for on other devices can be identified by analyzing the targeted advertising stored in seized devices. Using the proposed method, an investigator would extract the user interest while solving legal issues (e.g., privacy, jurisdiction) and technical issues (e.g., computing resources, the increase in storage capacity, accuracy) of traditional forensic techniques.

# Acknowledgement

# References

[1]  G. Palmer, "A road map for digital forensic research: report from the first digital forensic research workshop (DFRWS)," 2001 [Online]. Available: https://dfrws.org/wp-content/uploads/2019/06/2001_USA_a_road_map_for_digital_forensic_research.pdf

[2]  J. Hou, Y. Li, J. Yu, and W. Shi, "A survey on digital forensics in Internet of Things," *IEEE Internet of Things Journal*, vol. 7, no. 1, pp. 1-15, 2019.

[3]  A. Nieto and R. Rios, "Cybersecurity profiles based on human-centric IoT devices," *Human-centric Computing and Information Sciences*, vol. 9, article no. 39, 2019. https://doi.org/10.1186/s13673-019-0200-y

[4]  H. Arshad, A. B. Jantan, and O. I. Abiodun, "Digital forensics: review of issues in scientific validation of digital evidence," *Journal of Information Processing Systems*, vol. 14, no. 2, pp. 346-376, 2018.

[5]  L. Caviglione, S. Wendzel, and W. Mazurczyk, "The future of digital forensics: challenges and the road ahead," *IEEE Security & Privacy*, vol. 15, no. 6, pp. 12-17, 2017.

[6]  D. Jeong and S. Lee, "High-speed searching target data traces based on statistical sampling for digital forensics," *IEEE Access*, vol. 7, pp. 172264-172276, 2019.

[7]  J. I. James and P. Gladyshev, "A survey of mutual legal assistance involving digital evidence," *Digital Investigation*, vol. 18, pp. 23-32, 2016.

[8]  N. M. Karie and H. S. Venter, "Taxonomy of challenges for digital forensics," *Journal of Forensic Sciences*, vol. 60, no. 4, pp. 885-893, 2015.

[9]  A. Dehghantanha and K. Franke, "Privacy-respecting digital investigation," in *Proceedings of 2014 12th Annual International Conference on Privacy, Security and Trust*, Toronto, Canada, 2014, pp. 129-138.

[10]  Y. Xin, L. Kong, Z. Liu, Y. Chen, Y. Li, H. Zhu, M. Gao, H. Hou, and C. Wang, "Machine learning and deep learning methods for cybersecurity," *IEEE Access*, vol. 6, pp. 35365-35381, 2018.

[11]  Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436-444, 2015.

[12]  K. Chowdhary, "Natural language processing," in *Fundamentals of Artificial Intelligence*. New Delhi, India: Springer, 2020, pp. 603-649.

[13]  J. Hirschberg and C. D. Manning, "Advances in natural language processing," *Science*, vol. 349, no. 6245, pp. 261-266, 2015.

[14]  H. Kox, B. Straathof, and G. Zwart, "Targeted advertising, platform competition, and privacy," *Journal of Economics & Management Strategy*, vol. 26, no. 3, pp. 557-570, 2017.

[15]  S. Q. Liu and A. S. Mattila, "Airbnb: online targeted advertising, sense of power, and consumer decisions," *International Journal of Hospitality Management*, vol. 60, pp. 33-41, 2017.

[16]  J. R. C. Nurse and O. Buckley, "Behind the scenes: a cross-country study into third-party website referencing and the online advertising ecosystem," *Human-centric Computing and Information Sciences*, vol. 7, article no. 40, 2017. https://doi.org/10.1186/s13673-017-0121-6

[17]  M. Conti, V. Cozza, M. Petrocchi, and A. Spognardi, "TRAP: using targeted ads to unveil google personal profiles," in *Proceedings of 2015 IEEE International Workshop on Information Forensics and Security (WIFS)*, Rome, Italy, 2015, pp. 1-6.

[18]  S. Dhelim, N. Aung, and H. Ning, "Mining user interest based on personality-aware hybrid filtering in social networks," *Knowledge-Based Systems*, vol. 206, article no. 106227, 2020. https://doi.org/10.1016/j.knosys.2020.106227

[19]  J. Kang and H. Lee, "Modeling user interest in social media using news media and Wikipedia," *Information Systems*, vol. 65, pp. 52-64, 2017.

[20] X. Luo, J. Wang, Q. Shen, J. Wang, and Q. Qi, "User behavior analysis based on user interest by web log mining," in *Proceedings of 2017 27th International Telecommunication Networks and Applications Conference (ITNAC)*, Melbourne, Australia, 2017, pp. 1-5.

[21] N. N. Diep, N. Van Tien, N. H. Anh, and T. M. Phuong, "An unsupervised method for web user interest analysis," in *Proceedings of 2019 6th NAFOSTED Conference on Information and Computer Science (NICS)*, Hanoi, Vietnam, 2019, pp. 27-32.

[22] P. Siriaraya, Y. Yamaguchi, M. Morishita, Y. Inagaki, R. Nakamoto, J. Zhang, J. Aoi, and S. Nakajima, "Using categorized web browsing history to estimate the user's latent interests for web advertisement recommendation," in *Proceedings of 2017 IEEE International Conference on Big Data (Big Data)*, Boston, MA, 2017, pp. 4429-4434.

[23] E. B. Karbab and M. Debbabi, "MalDy: portable, data-driven malware detection using natural language processing and machine learning techniques on behavioral analysis reports," *Digital Investigation*, vol. 28, pp. S77-S87, 2019. https://doi.org/10.1016/j.diin.2019.01.017

[24] D. H. Lee, Y. R. Kim, H. J. Kim, S. M. Park, and Y. J. Yang, "Fake news detection using deep learning," *Journal of Information Processing Systems*, vol. 15, no. 5, pp. 1119-1130, 2019.

[25] J. Salminen, M. Hopf, S. A. Chowdhury, S. G. Jung, H. Almerekhi, and B. J. Jansen, "Developing an online hate classifier for multiple social media platforms," *Human-centric Computing and Information Sciences*, vol. 10, article no. 1, 2020. https://doi.org/10.1186/s13673-019-0205-6

[26] E. Akbal, F. Gunes, and A. Akbal, "Digital forensic analyses of web browser records," *Journal of Software*, vol. 11, no. 7, pp. 631-637, 2016.

[27] C. Flowers, A. Mansour, and H. M. Al-Khateeb, "Web browser artefacts in private and portable modes: a forensic investigation," *International Journal of Electronic Security and Digital Forensics*, vol. 8, no. 2, pp. 99-117, 2016.

[28] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, "Scalable object detection using deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, Columbus, OH, 2014, pp. 2155-2162.

[29] Facebook Research, "fastText," 2017 [Online]. Available: https://github.com/facebookresearch/fastText/

[30] Keras team, "Keras," 2021 [Online]. Available: https://github.com/keras-team/keras

[31] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, vol. 25, pp. 1097-1105, 2012.

**Wonkyung Kim** https://orcid.org/0000-0003-1892-8600

He received the B.S. degree from the Division of Computer Science, Korea University in 2020. He is currently an M.S. student in the Institute of Cyber Security & Privacy at the Korea University (Korea). His research interests include digital forensics, artificial intelligence, and digital profiling.

**Kukheon Lee** https://orcid.org/0000-0002-1512-4968

He is a Ph.D. student at the Institute of Cyber Security & Privacy at Korea University. His research revolves mostly around digital forensics (in particular, multimedia and database). He has a lot of project experience. He devised forensic accounting method and tool for account data in a database (with Supreme Prosecutors' Office of Korea) and recovery system for the corruption of document files (with Police of Korea). Also, he developed a recovery tool for the embedded video systems.

**Sangjin Lee** https://orcid.org/0000-0002-6809-5179

He received a Ph.D. degree from the Department of Mathematics, Korea University, in 1994. From 1989 to 1999, he was a Senior Researcher with the Electronics and Telecommunications Research Institute, South Korea. He has been running the Digital Forensic Research Center, Korea University, since 2008. He is currently the President of the Division of Information Security, Korea University. He has authored or coauthored over 130 articles in various archival journals and conference proceedings and over 200 articles in domestic journals. His research interests include digital forensics, data processing, forensic framework, incident response, and so on.

**Doowon Jeong** https://orcid.org/0000-0002-2997-2335

He received the B.S. degree from the Division of Industrial Management Engineering, Korea University, in 2011, and the Ph.D. degree from the Graduate School of Information Security, Korea University, in 2019. He is currently an assistant professor with the College of Police and Criminal Justice, Dongguk University. His research interests include digital forensics, information security, artificial intelligence, and digital profiling.