

PAPER • OPEN ACCESS

## Analysis WhatsApp Forensic and Visualization in Android Smartphone with Support Vector Machine (SVM) Method

To cite this article: Ubaidillah *et al* 2019 *J. Phys.: Conf. Ser.* **1196** 012064

View the [article online](#) for updates and enhancements.



**IOP | ebooks™**

Bringing you innovative digital publishing with leading voices to create your essential collection of books in STEM research.

Start exploring the **collection** - download the first chapter of every title for free.

# Analysis WhatsApp Forensic and Visualization in Android Smartphone with Support Vector Machine (SVM) Method

Ubaidillah<sup>1</sup>, D Stiawan<sup>1</sup>, T W Septian<sup>1</sup>, R F Malik<sup>1</sup>, Fepiliana<sup>1</sup>, A Heriyanto<sup>1</sup>,  
R Budiarto<sup>2</sup>

<sup>1</sup> The Faculty of Computer Science, University of Sriwijaya, Palembang, ID

<sup>2</sup> The Faculty of Computer Information System, Al Baha University, Al Baha, SA

Email: ubai.com@gmail.com

**Abstract.** In this paper we discuss about the pattern of data WhatsApp application in the form of text, image, and sound. This application is widely used as a media to deliver and exchange information, so this is also open opportunity to commit crimes through this media. In this research, we focus on forensic and create visual about patterns data from WhatsApp in form of text, image, audio and video. In this study, 288 instances have been examined using 9 features. The feature that characterizes any personal data in the selected data has not been used. From the results, we get by a reading database which stored on the client smartphone is obtained results of the conversation time, client contact number, key and message types between client applications use WhatsApp. Receiver Operating Characteristic (ROC) curves was also presented for each of the classes.

## 1. Introduction

WhatsApp became one of the most used personal mobile messaging as delivery of text and free content (which are audio, video, pictures, location, and contacts) [1]. The effect of high usage of WhatsApp application caused due to the large community which makes it as a platform for them to communicate. On the other hand [2], a short message is one of the important components of proving evidence in many cases (which high-profile). Various kinds of criminal investigation and trial are very helped with the digital evidence from a smartphone.

The artifacts were left in WhatsApp message can be analyzed by forensics. Then it will give a complete description of all the artifacts which generate by those WhatsApp messages, we'll talk about decoding and interpreters. We will obtain various types of information by correlated one by one, which could not be gained from their isolation system.

Short messaging application service is very potential to get the source of evidence information for most investigations. Because of the short message service is used not only for good things but also for the crime. The crime interacts with victims also to escape and avoid arrest [3].

The methods and tools which used in this research are highly relevant to the investigation of the case, to prove a call or message content is done on a particular date and time can be used as a series of stories in a case. For this research will be presented forensic steps which conducted such as recognizing the artifact, artifact extraction, decryption, and analysis of data from WhatsApp application. Then we use SVM to classify the data obtained from extracting log files and exported databases from the WhatsApp application.



## 2. Relevant theory

Much forensic research which conducted on WhatsApp, but most focus on the data which is stored on the mobile device but does not display the data process in visualization on how such data can be interpreted and decoding. In this section, we will present some previous research literature review.

### 2.1. Mobile Device Forensic

In the explanation [4] and [5], both focus on WhatsApp analysis forensic on android. This study found the forensic acquisition of the artifacts which left by WhatsApp on the device. In [4], research focuses on forensic analysis about storage on internal memory or external. The results show that someone was able to gain a lot of artifacts such as phone numbers, messages, media files, location, picture profiles, log and more. While [5] is analyzing artifacts WhatsApp and Viber using toolkit Forklatic Forensic Extraction Device (UFED). They can restore the list of contacts, messages and media exchanges.

On the other hand, [6] has done analysis against WhatsApp on Android devices. The paper gives a artifacts description which comprehensive which generated by WhatsApp and discusses parsing, interpretation and the relationships between artifacts. So it is able to provide an analysis on how to reconstruct the contact list and messages chronology which have been exchanged by the user.

On [7] have checked 20 popular mobile social messaging application for android (one of which namely WhatsApp). In their research focus on the traffic which is not encrypted, so that it can be easily reconstructed. On the application network traffic encryption, WhatsApp more favorable compared to other social messaging application.

Further research [8][9] focus on the backup analysis of iTunes iOS devices aims to identify artifacts left by various social networking applications, including WhatsApp. Their research only analyzes databases chat WhatsApp and only a portion of the chat which was analyzed.

WhatsApp has a feature to display a notification when a message is received and read by the recipient, it is useful to the investigation because it would give evidence to someone who claimed that they did not read a message [10]. WhatsApp also has the ability to transfer a wide range of media files, such as images, audio, and video. The chances of recovery media files sent and received can be investigated, it's also useful to see who and when the sending and receiving files.

### 2.2. Support Vector Machine (SVM)

Many machine learning algorithms can be used to classify types of the fragment, such as support vector machine (SVM), k-nearest neighbors (k-NN), neural networks (ANN) and so on. On Research [11] that the k-NN and ANN algorithm found not as effective as the SVM algorithm. In addition, the SVM algorithm has also been widely used by previous researchers.

Classification task usually involves the separation of the data into training and testing sets. Each example in the training group contained one "target value" (label class) and some "attributes" (fragment features). The purpose of SVM is generating models (based on the training data) which predict the target value of testing data which is only being given the attribute testing data.

There are four basic kernels of SVM:

$$\text{Linear:} \quad K(x_i, x_j) = x_i^T x_j \quad (1)$$

$$\text{Polynomial:} \quad K(x_i, x_j) = (\gamma x_i^T x_j + r)^d, \gamma > 0 \quad (2)$$

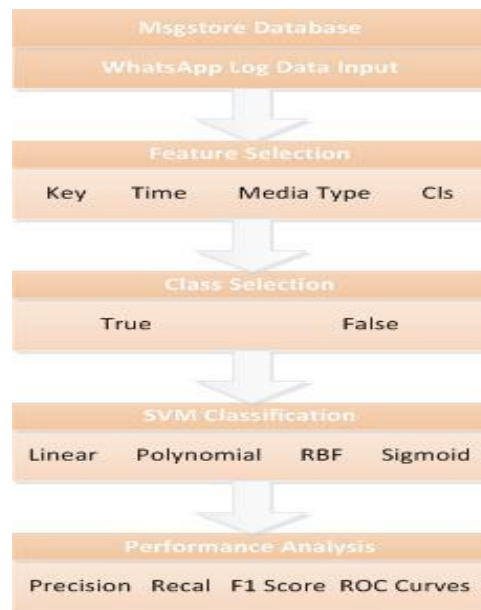
$$\text{Radial Basis Function (RBF):} \quad K(x_i, x_j) = \exp(-\gamma \|x_i - x_j\|^2), \gamma > 0 \quad (3)$$

$$\text{Sigmoid:} \quad K(x_i, x_j) = \tanh(\gamma x_i^T x_j + r) \quad (4)$$

Where  $\gamma$   $r$  and  $d$  is kernel parameter. Need to select a kernel function while using SVM. Normally function Linear kernel and radial used first. In [12] found that the linear kernel is more effective than a radial basis function (RBF). SVM has a generalization and very good scaling properties [13]. Classification types of files and data can help forensic test and improve the efficiency of the investigation, with the help investigators target or prioritize search efforts, extract, and analysis on the types of bigger interests toward the case. The case can be either type-criminal, civil, regulatory, or simply a matter of Organization (e.g., harassment, fraud, abuse, etc.).

### 3. Research method

In this paper we use software-emulated Android devices in place of physical ones. In particular, we use the NoxPlayer [14] that is able to faithfully emulate the behavior of a complete Android device. NoxPlayer implements the internal device memory as a Virtual Box storage file, whose format is documented and, therefore, can be parsed by a suitable tool to extract the files stored inside it. In this way, the acquisition of the internal memory of the device is greatly simplified, as it reduces to inspect the content of this file.



**Figure 1.** Proposed Method

Our tool is a command-line program written in Python (version 2.7). It converts the log file exported from NoxPlayer to a CSV file. It is available in the form of source code. It requires one input parameter to a CSV file containing the details of dissected packets. Table 1 shows the descriptions of the selected attributes.

**Table 1.** Features and Descriptions

Feature	Descriptions
timestamp	Transaction time session ended
key	Sender number and Receiver number
from_me	TRUE and FALSE (true = sender; false = receiver)
media_wa_type	Type of media in WhatsApp
remote_jid	WhatsApp ID of the communication partner
cls	Message type has been sending and receiving
status	Account state on contact application
file_type	Files type
has_video	Messages that have videos or without videos

There are 2 classes in the from me attribute used as a class. Descriptions of these classes are shown in Table 2.

**Table 2.** Activities send and receive data WhatsApp

From me	Description
TRUE	Message from the sender
FALSE	Message from the receiver

#### 4. Result and discussion

In an experiment carried out using WhatsApp application with normal conditions which are without modifying the application usage, the meaning of the first data obtained is suitable with the conditions of WhatsApp application use so we take all the data files as a whole or complete for either data conversation, sending pictures or conversations (voice and video call).

From the stages which we do above all is used as a media to get the raw data. The results of the raw data which we will use to do extraction to read patterns of messages which sent either in the form of text, sound, images and video (see figure 2).

```
2018-10-11 17:06:08.550 LL_I W [158:Signal Protocol] axolotl received a message; message.key=Key
[id=890FB5FE43F81046472B3B096981C91E, from_me=false, remote_jid=6285369933878@s.whatsapp.net];
message.retryCount=0; message.remote_resource=
2018-10-11 17:06:08.572 LL_I W [158:Signal Protocol] msgstore/add/recv; key=Key
[id=890FB5FE43F81046472B3B096981C91E, from_me=false, remote_jid=6285369933878@s.whatsapp.net];
media_wa_type=2; status=0
2018-10-11 17:06:08.589 LL_I W [158:Signal Protocol] msgstore/unseen/1/1/0/1539252368628
2018-10-11 17:06:08.617 LL_I W [158:Signal Protocol] axolotl updating session for 6285369933878
2018-10-11 17:06:08.641 LL_I W [158:Signal Protocol] axolotl stored session for 6285369933878
2018-10-11 17:06:11.500 LL_I W [157:ReaderThread] xmp/reader/read/status-update-from-target Key
[id=D0DF8BB054EA3D366EF6AEA83299D420, from_me=true, remote_jid=6285768044255-1539242223@g.us]
```

Figure 2. Events in the log file corresponding to (a) from me, (b) media\_wa\_type, (c) status.

From the file\_type attribute, there are 8 categories (audio, contact, image, location, null, ptt, none, video) and from the has\_video attribute there are 3 numeric (0,1,2) also from the media\_wa\_type attribute obtained by 7 numerically (0,1,2, 4,5,8,10) and from the status attribute 3 numerical values (0,1,6) are obtained. The graphs of the obtained extraction raw data as shown in figure 3.

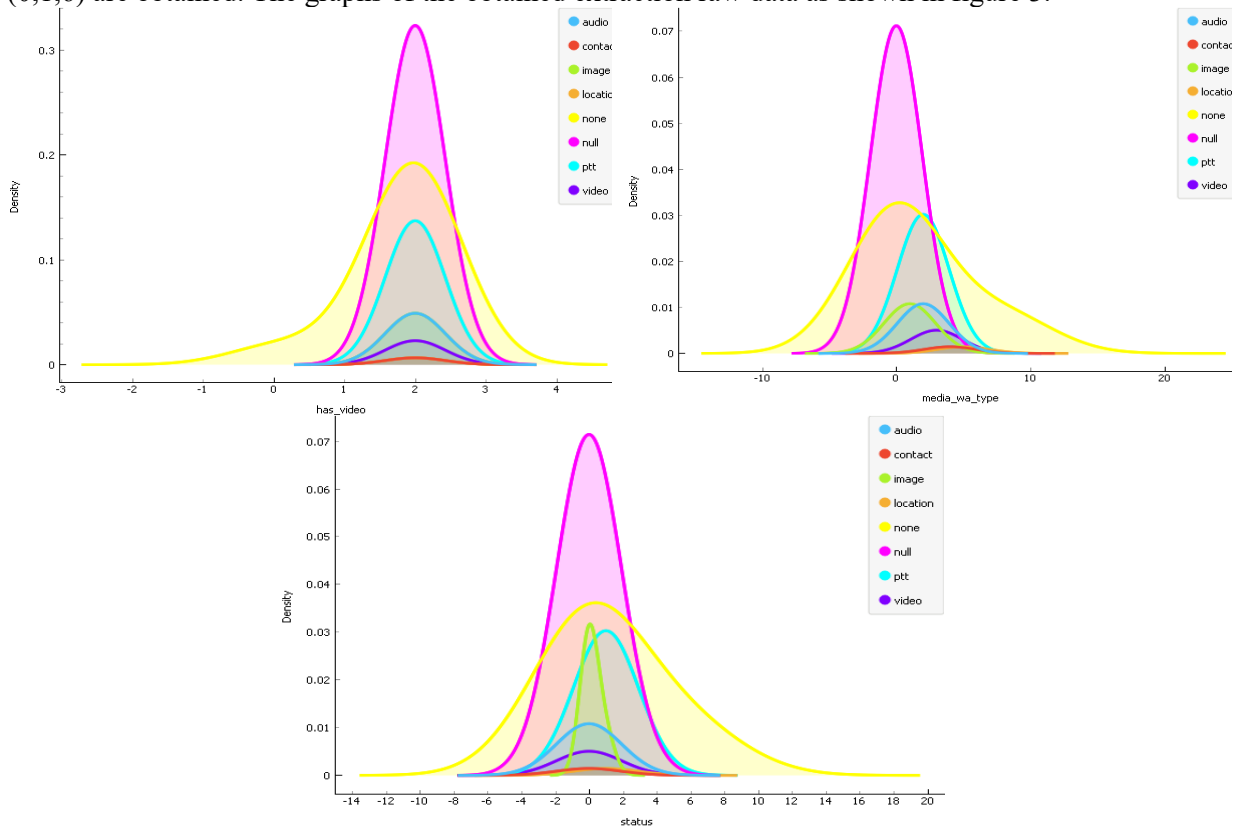


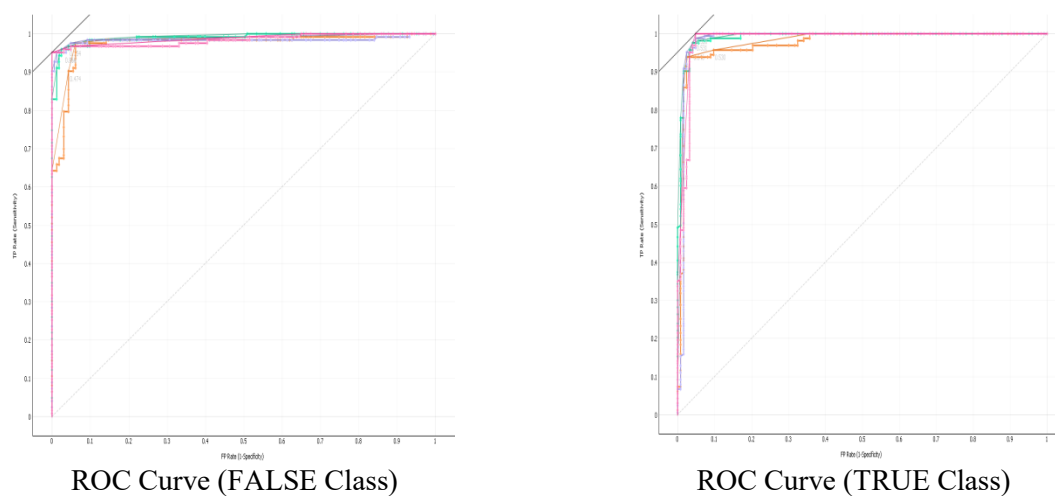
Figure 3. Visualization of attribute has\_video, media\_wa\_type, and status

For the  $F_1$  score, it is seen that the worst result is obtained with the classifier selected as the polynomial activation function as in the precision and recall values.

**Table 3.** Evaluation Results

Method	F1	Precision	Recall
SVM Sigmoid	97.6	97.7	97.6
SVM RBF	97.2	97.2	97.2
SVM Polynomial	91.6	91.8	91.6
SVM Linear	95.8	95.8	95.8

The graphs of the obtained ROC curves are shown in figure 4.



**Figure 4.** The graphs of the obtained ROC curves

## 5. Summary

In this research, we focus on application condition WhatsApp with two condition text and voice with the purpose to look and generating attributes which used in sending process data. On the data which we get to the number of attributes for text obtained eight attributes, while the voice data obtained attribute as much as sixteen attributes. Further, in the future, we will develop a classification system using machine learning method so that the data which obtained are better. The study was conducted using 288 instances and 9 features. The action, which is one of the attributes used, is selected as the class. The values of this class are “true”, “false”. It was observed that the highest recall value was obtained in the SVM classifier with 97.7% of the Sigmoid activation function selected and the highest precision value was obtained in the SVM classifier with the Sigmoid activation function of 97.6%. As the  $F_1$  score value, it was observed that the best result was achieved with the classifier in which Sigmoid activation was used with 97.6%.

## References

- [1] Nicole Arce, “WhatsApp Calling For Android And iOS: How To Get It And What To Know,” 2015. [Online]. Available: <http://www.techtimes.com/articles/38291/20150309/whatsapp-calling-for-android-and-ios-how-to-get-it-and-what-to-know.htm>.
- [2] D. Walnycky, I. Baggili, A. Marrington, J. Moore, and F. Breitingner, “Network and device forensic analysis of Android social-messaging applications,” *Digit. Investig.*, vol. 14, no. S1,

- pp. S77–S84, 2015.
- [3] C. Anglano, M. Canonico, and M. Guazzone, “Forensic analysis of Telegram Messenger on Android devices,” *Digit. Investig.*, pp. 1–7, 2017.
  - [4] N. S. Thakur, “Forensic Analysis of WhatsApp on Android Smartphones,” 2013.
  - [5] A. Mahajan, M. S. Dahiya, and H. P. Sanghvi, “Forensic Analysis of Instant Messenger Applications on Android Devices,” *Int. J. Comput. Appl.*, vol. 68, no. 8, pp. 38–44, 2013.
  - [6] F. Karpisek, I. Baggili, and F. Breitter, “WhatsApp network forensics: Decrypting and understanding the WhatsApp call signaling messages,” *Digit. Investig.*, vol. 15, pp. 110–118, 2015.
  - [7] M. I. Husain and R. Sridhar, “iForensics: Forensic analysis of instant messaging on smart phones,” *Lect. Notes Inst. Comput. Sci. Soc. Telecommun. Eng.*, vol. 31 LNICST, no. Vim, pp. 9–18, 2010.
  - [8] Y. Tso, S.-J. Wang, C.-T. Huang, and W. Wang, “iPhone social networking for evidence investigations using iTunes forensics,” *Proc. 6th Int. Conf. Ubiquitous Inf. Manag. Commun. - ICUIMC '12*, p. 1, 2012.
  - [9] N. B. Al Barghuthi and H. Said, “Social networks IM forensics: Encryption analysis,” *J. Commun.*, vol. 8, no. 11, pp. 708–715, 2013.
  - [10] A. Shortall and M. A. H. Bin Azhar, “Forensic Acquisitions of WhatsApp Data on Popular Mobile Platforms,” *Proc. - 2015 6th Int. Conf. Emerg. Secur. Technol. EST 2015*, pp. 13–17, 2016.
  - [11] N. Zheng, T. Wu, J. Wang, and M. Xu, “A fragment classification method depending on data type,” *Proc. - 15th IEEE Int. Conf. Comput. Inf. Technol. CIT 2015, 14th IEEE Int. Conf. Ubiquitous Comput. Commun. IUCC 2015, 13th IEEE Int. Conf. Dependable, Auton. Se.*, pp. 1948–1953, 2015.
  - [12] S. Fitzgerald, G. Mathews, C. Morris, and O. Zhulyn, “Using NLP techniques for file fragment classification,” *Digit. Investig.*, vol. 9, no. SUPPL., 2012.
  - [13] N. L. Beebe, L. A. Maddox, L. Liu, and M. Sun, “Sceadan: Using concatenated N-gram vectors for improved file and data type classification,” *IEEE Trans. Inf. Forensics Secur.*, vol. 8, no. 9, pp. 1519–1530, 2013.
  - [14] “Mainkan game Android di PC ataupun Mac dengan menggunakan Emulator Android NoxPlayer.” [Online]. Available: <https://id.bignox.com/>. [Accessed: 23-Jul-2018].