

# Fuzzy based binary feature profiling for modus operandi analysis

Mahawaga Arachchige Pathum Chamikara<sup>1,2</sup>, Akalanka Galappaththi<sup>1</sup>,  
Roshan Dharshana Yapa<sup>1,2</sup>, Ruwan Dharshana Nawarathna<sup>1,2</sup>,  
Saluka Ranasinghe Kodituwakku<sup>1,2</sup>, Jagath Gunatilake<sup>1,2</sup>,  
Aththanapola Arachchilage Chathranee Anumitha Jayathilake<sup>2</sup> and  
Liwan Liyanage<sup>3</sup>

<sup>1</sup> Postgraduate Institute of Science (PGIS), University of Peradeniya, Peradeniya, Sri Lanka

<sup>2</sup> Faculty of Science, University of Peradeniya, Peradeniya, Sri Lanka

<sup>3</sup> School of Computing, Engineering and Mathematics, University of Western Sydney, Western Sydney, NSW, Australia

## ABSTRACT

It is a well-known fact that some criminals follow perpetual methods of operations known as modi operandi. Modus operandi is a commonly used term to describe the habits in committing crimes. These modi operandi are used in relating criminals to crimes for which the suspects have not yet been recognized. This paper presents the design, implementation and evaluation of a new method to find connections between crimes and criminals using modi operandi. The method involves generating a feature matrix for a particular criminal based on the flow of events of his/her previous convictions. Then, based on the feature matrix, two representative modi operandi are generated: complete modus operandi and dynamic modus operandi. These two representative modi operandi are compared with the flow of events of the crime at hand, in order to generate two other outputs: completeness probability (CP) and deviation probability (DP). CP and DP are used as inputs to a fuzzy inference system to generate a score which is used in providing a measurement for the similarity between the suspect and the crime at hand. The method was evaluated using actual crime data and ten other open data sets. In addition, comparison with nine other classification algorithms showed that the proposed method performs competitively with other related methods proving that the performance of the new method is at an acceptable level.

Submitted 21 November 2015

Accepted 12 May 2016

Published 13 June 2016

Corresponding author

Mahawaga Arachchige  
Pathum Chamikara,  
pathumchamikara@gmail.com

Academic editor

Sebastian Ventura

Additional Information and  
Declarations can be found on  
page 27

DOI 10.7717/peerj-cs.65

© Copyright  
2016 Chamikara et al.

Distributed under  
Creative Commons CC-BY 4.0

OPEN ACCESS

**Subjects** Algorithms and Analysis of Algorithms, Artificial Intelligence, Data Mining and Machine Learning

**Keywords** Modus operandi analysis, Fuzzy inference systems, Binary feature analysis, Classification, Association rule mining

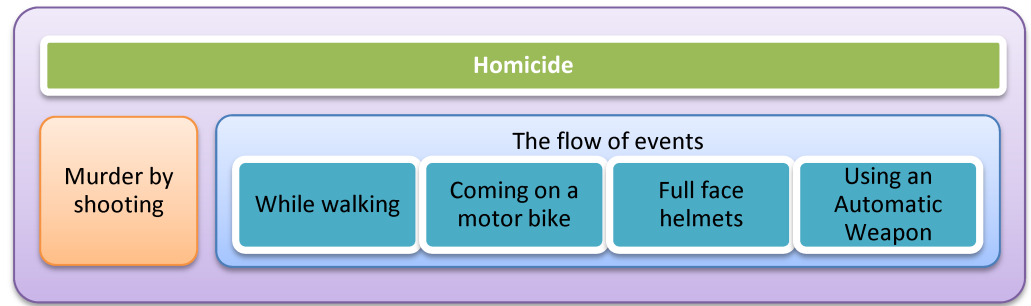
## INTRODUCTION

Scientists have long played a role in examining deviant behavior in society. “Deviance behaviour” is a term used by scientists to refer to some form of “rule-breaking” behaviour (*Holdaway, 1993*). It can be the behaviour of violating a social norm or the law. Criminal behaviour is also a form of deviance, one that is defined as the breaking of legal rules. Nevertheless, there is a difference between deviance and crime. Deviance involves breaking a norm and evoking a negative reaction from others. Crime is a deviance that

breaks a law, which is a norm stipulated and enforced by government bodies ([Holdaway, 1993](#)). However, crimes negatively affect society. Therefore, law enforcement authorities take necessary actions to mitigate crimes in an environment where high crime frequencies are observed each year. In this exercise, the application of technology for crime analysis is being widened in the world. Locard's exchange principle states that every contact of the perpetrators of a crime scene leaves a trace. The perpetrators will both bring something into the scene and leave with something from the scene ([Chisum & Turvey, 2000](#)). However, the cognitive abilities of criminals will always make them minimize their risks of apprehension by conducting the perfect crime and maximizing their gain ([Paternoster & Bachman, 2001](#)). Modus operandi or method of operation such as preparation actions, crime methods and weapons are frequently used in criminal profiling because the past crime trends show that, after criminals get used to a certain method of operation, they try to use the same modus operandi in committing his/her next crime ([Palmiotto, 1988](#)).

The criminals develop a set of actions during the performance of a series of crimes which we refer to as "modus operandi" (MO). MO is developed with the crimes he/she commits and the nature of trying to stick with the developed MO that has worked throughout the previous crimes ([Douglas & Douglas, 2006](#)). In any criminal career, the MO happens to evolve, no matter what the circumstances. Also, it is a common behaviour that serial offenders tend to exhibit significant behaviour known as his/her signature. Therefore, MOs of criminals play a major role in investigating crimes ([Douglas & Douglas, 2006](#)). It is a known fact that features such as criminal signature and physical appearance are used in crime investigations in almost all the police departments around the world. Sri Lanka police also use MOs of criminals to identify the suspects who have conducted crimes. Currently Sri Lanka Police use a manual crime recording and investigation system. This manual system has many problems such as data redundancy, inefficiency, tediousness, inability to support crime investigation and many other problems which are associated with a conventional manual system. To overcome these problems, a web-based framework was proposed with geographical information support containing a centralized database for crime data storage and retrieval, named SL-CIDSS: Sri Lanka Crime Investigation Decision Support System ([Chamikara et al., 2015](#)). The proposed system accompanies a collection of data mining algorithms which effectively support the crime investigation process. Fuzzy based binary feature profiling (BFPM) for modus operandi analysis is one novel algorithm which is integrated with the system to provide an effective way to find the similarity between crimes and criminals.

According to the penal code of Sri Lanka first enacted in 1882 and amended subsequently several times in later years ([The 'Lectric Law Library, 0000](#)), Sri Lanka police classifies crimes into two categories: Grave crimes and Minor offences. Until 2014, grave crimes were classified under 21 crime categories and in 2015 another 5 new crime categories were introduced, making it 26 categories of grave crime types. Kidnapping, Fraud or mischief causing damage greater than 25,000 rupees, Burglary, Grievous hurt, Hurt by sharp weapon, Homicide, Rape, Robbery, Cheating by trust, Theft are 10 of the most frequent crime types. To identify the patterns involved in crimes, a collection of subtypes were identified under these 26 crime types. These subtypes have been created mainly for the purpose of modus



**Figure 1** Relationship between main crime type, subtypes and crime flows.

operandi analysis. Most frequent behaviors of criminals/crimes are considered as crime subtypes. When a crime is logged in the Grave Crime Record (GCR) book, it is classified under one of the 26 main categories. But, under the section of “nature of crime” in the GCR book, the police officers record the flow of the crime incident including the subtypes.

A subtype is a sub category of one of the main crime types. For investigation, the nature of the crime is broken into subtypes and flows according to their frequency of occurrence and uniqueness. These sub categorizations have been introduced mainly to minimize the broadness of main type and to improve clarity. [Figure 1](#) depicts the relationship of the subtypes and flows where there can be a flow of events to a crime recorded as one of the 26 main crime types. For the simplicity and easy handling of data, the investigators have provided subtype codes and flow codes. The flow of events provides a modus operandi which is most of the time unique to an offender. Each subtype is provided with a code under the main type, to make the crime investigation process easier. For example, ROB/S001 denotes a subtype that is Highway robbery; here ROB denotes the main type under which the corresponding subtype appears. In this case, it is Robbery. Crime types are further subdivided into sub types to make the analysis and processing simpler. In this manner, crime subtypes and flows have been identified under all the 26 crime types. The space for adding more subtypes and flows under these crime types exists. A new subtype or a flow is introduced to a particular main crime, if the same subtype or the flow happens to persist for a prolonged time.

This paper proposes a novel method of criminal profiling using modus operandi which can be used to identify associations between crimes and criminals. The method is based on a new technique named, “binary feature vector profiling.” Key relationships between a criminal and the previous convictions are analyzed using binary feature profiling and association rule mining techniques. Due to the impreciseness and vagueness of these extracted attributes, a fuzzy inference system is used in making the final decision. The newly proposed method was adapted into a classification algorithm in order to test its accuracy. An actual crime data set was used in testing the performance of the newly proposed method and it was compared against nine well-established classification algorithms using ten open data sets. The results confirmed that the proposed method produce competitive results compared to the other nine classification algorithms.

The rest of the paper is organized as follows. The Related work section presents a summary of the work that has been conducted on modus operandi analysis as well as a brief discussion on crime investigation using link analysis and association mining in general. The Materials and Methods section discusses the main steps of the newly proposed algorithm. Next, the Results and Discussion section provides a validation and performance evaluation of the newly proposed method along with a performance comparison with nine other classification algorithms. Finally, some concluding remarks and future enhancements are outlined in the Conclusion section.

## Related work

Literature shows many methods which have been developed in the area of automated crime investigation. Our major concern has been laid upon the research carried out on crime investigation using association mining as our research considers on developing a model to find the associations between the criminals and the crimes depending on the modes operandi. *Bennell & Canter (2002)* have proposed a method to use statistical models to test directly the police practice of utilizing modus operandi to link crimes to a common offender. The results indicated that certain features such as the distance between burglary locations, lead to high levels of predictive accuracy. *Bennell & Jones (2005)* have tried to determine if readily available information about commercial and residential serial burglaries, in the form of the offender's modus operandi, provides a statistically significant basis for accurately linking crimes committed by the same offenders. *Leclerc, Proulx & Beauregard (2009)* have reviewed the theoretical, empirical, and practical implications related to the modus operandi of sexual offenders against children. They have presented the rational choice perspective in criminology followed by descriptive studies aimed specifically at providing information on modus operandi of sexual offenders against children.

Clustering crimes, finding links between crimes, profiling offenders and criminal network detection are some of the common areas where data mining is applied in crime analysis (*Oatley & Ewwart, 2011; King & Sutton, 2013; Borg et al., 2014*). Association analysis, classification and prediction, cluster analysis, and outlier analysis are some of the traditional data mining techniques which can be used to identify patterns in structured data. Offender profiling is a methodology which is used in profiling unknown criminals or offenders. The purpose of offender profiling is to identify the socio-demographic characteristics of an offender based on information available at the crime scene (*Mokros & Alison, 2002; Canter et al., 2013*). Association rule mining discovers the items in databases which occur frequently and present them as rules. Since this method is often used in market basket analysis to find which products are bought with what other products, it can also be used to find associated crimes conducted with what other crimes. Here, the rules are mainly evaluated by the two probability measures, support and confidence (*Agrawal, Imielinski & Swami, 1993; Yi et al., 2015*). Association rule mining can also be used to identify the environmental factors that affect crimes using the geographical references (*Koperski & Han, 1995*). Incident association mining and entity association mining are two applications of association rule mining. Incident association mining can be used to find the crimes committed by the same offender and then the unresolved crimes can be linked to find the

offender who committed them. Therefore, this technique is normally used to solve serial crimes like serial sexual offenses and serial homicides (Chen, 2006).

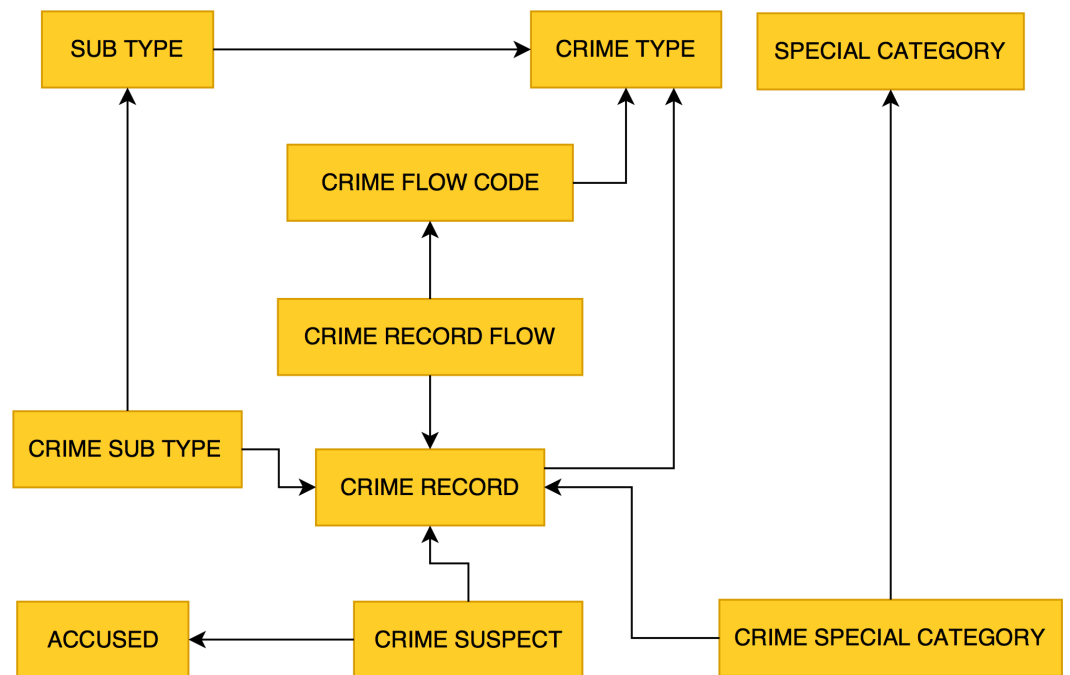
Similarity-based association mining and outlier-based association mining are two approaches used in incident association mining. Similarity-based association mining is used mainly to compare the features of a crime with the criminal's behavioral patterns which are referred as modus operandi or behavioral signature. In outlier-based association mining, crime associations will be created on the fact that both the crime and the criminal have the possibility of having some distinctive feature or a deviant behavior (Lin & Brown, 2006). Entity association mining/link analysis is the task of finding and charting associations between crime entities such as persons, weapons, and organizations. The purpose of this technique is to find out how crime entities that appear to be unrelated at the surface, are actually linked to each other (Chen, 2006). Link analysis is also used as one of the most applicable methods in social network analysis (Berry & Linoff, 2011) in finding crime groups, gate keepers and leaders (Chen et al., 2003).

Attribution can be used to link crimes to offenders. If two offences in different places involve the same specific type, those may be readily attributed to the same offender (Oatley & Ewwart, 2011). There are three types of link analysis approaches, namely Heuristic-based, Statistical-based and Template-based (Chen, 2006). Sequential pattern mining is also a similar technique to association rule mining. This method discovers frequently occurring items from a set of transactions occurred at different times (Chen et al., 2014). Deviation detection detects data that deviates significantly from the rest of the data which is analyzed. This is also called outlier detection, and is used in fraud detection (Chen et al., 2014; Capozzoli, Lauro & Khan, 2015).

In classification, the data points will be assigned to a set of predefined classes of data by identifying a set of common properties among them. This technique is often used to predict crime trends. Classification needs a reasonably complete set of training and testing data since a high degree of missing data would limit the prediction accuracy (Chen et al., 2014). Classification comes under supervised learning method (Chen, 2006; Chikersal, Poria & Cambria, 2015) which includes methods such as Bayesian models, decision trees, artificial neural networks (Chen, 1995) and support vector machines. String comparison techniques are used to detect the similarity between the records. Classification algorithms compare the database record pairs and determine the similarity among them. This concept can be used to avoid deceptive offender profiles. Information of offenders such as name, address, etc. might be deceptive and therefore the crime database might contain multiple records of the same offender. This makes the process of identification of their true identity difficult (Chen et al., 2014).

## SYSTEMS AND METHODS

This section provides a description about the systems and methods used in developing the fuzzy based binary feature profiling for modus operandi analysis. First, an overview about how SL-CIDSS captures the logics of modus operandi is explained. Then a detailed description about the steps of the newly proposed algorithm is explained.



**Figure 2** Crime flow entity arrangement in SL-CIDSS.

Figure 2 shows how SL-CIDSS database captures the crime types and subtypes. A crime record has a crime record flow. Typically, a crime is committed by a criminal and a particular accused might commit one or more crimes. A CRIME RECORD can be of one the 26 crime types. A particular CRIME RECORD will be considered under one main CRIME TYPE with the highest precedence in the order of seriousness. For example, a crime incident that includes a murder and a robbery will be categorized as a murder though a robbery has also taken place. But in the nature of crime section, all crimes followed by the main type will be stated. Therefore, the CRIME RECORD FLOW captures all the steps of the crime as a sequence of steps recorded. The crime flows that have been previously registered are mapped under CRIME FLOW CODE. Also, a particular CRIME RECORD instance can contain multiple SUB TYPES which are recorded as CRIME SUB TYPE. The SPECIAL CATEGORY captures the crimes with special features such as crimes occurring at the same location or retail shop. A crime may involve several special categories which are saved in the CRIME SPECIAL CATEGORY. The ACCUSED entity records the information of suspects and accused and they are related to crime through the CRIME SUSPECT entity.

As the first step of the newly employed method, a feature matrix is generated, resulting in a binary matrix representing the crime flows. This binary feature matrix is composed of the binary patterns generated on previous convictions of a particular criminal/suspect. This binary form of the feature matrix provides a provision to direct application of computer algorithms with methods such as Apriori based association rule mining. The reduced complexity of the binary feature matrices provides an easy manipulation over the categorical and continuous valued features. Figure 3 shows the steps of the proposed MO analysis algorithm.



- Step 1:** Generate the feature matrix.
- Step 2:** Generate the dynamic MOs (DMO) of the criminals.
- Step 3:** Generate the complete MO profile (CMOP) of the criminals.
- Step 4:** Find the deviation probability (DP) of CMOP from the crime MO under consideration (UMO).
- Step 5:** Find the completeness probability of UMO against DMO.
- Step 6:** Use the two values obtained from step 4 and 5 as inputs of a fuzzy Inference system to obtain the final similarity value (out of 100).
- Step 7:** Classify the UMO under the class with highest similarity score for validation.

**Figure 3** Steps of the newly employed algorithm.

### Generating the feature matrix

[Table 1](#) shows how the feature vectors are generated and provides the way to generate modi operandi of criminals as binary sequences. According to the table, events of the crime scene are observed starting from its crime type. After a particular crime type is identified, the feature vectors are updated with ones for each subtype and flow code that is available in the crime or suspect's modus operandi. The vectors will be filled by zeros in places which the modus operandi does not have any contact with. The column names to the feature matrix are generated in such a way that it covers the collection of main types, sub types, crime flows and special categories at hand. For example, if we consider the list of crime types, subtypes, crime flows and the special category in [Table 1](#), it results in 21-bit feature vectors as shown in the last two columns.

In this manner we can produce binary MO patterns based on the crimes committed by different criminals as shown in the last two columns of [Table 1](#). According to [Table 1](#), Suspect 1 has committed a robbery with the subtypes, ABD/S003 (an abduction of a child from the legal guardian), ROB/S001 (an organized vehicle robbery) and the flows, ROB/F001 (Identity cards have been shown), ROB/F003 (accused has been wearing uniforms). Suspect 2 has committed a house breaking with the sub type BGL/S004 (use of stealth), and the flows, BGL/F001 (Entering from the window), BGL/F003 (Removing Grills).

[Table 2](#) shows a feature matrix of binary patterns which is generated by considering the previous convictions of suspect 1 assuming that he has conducted another robbery (conviction 2). *ct*, *st*, *fl* and *sc* in [Table 2](#) represent the abbreviations for "crime type," "sub type," "crime flow" and "special category" respectively.

### Generating the dynamic MOs (DMOs) of the criminals

Dynamic MO is a binary feature vector which is generated on bit patterns of the feature matrix of a particular criminal. The main purpose of the DMO is to obtain a criminal specific crime flow which captures the crime patterns which are frequently followed by a particular criminal. It is named as the dynamic modus operandi as it is subject to change when the new crime flows are added to the feature matrix. Therefore, this addresses the changing nature

**Table 1** An instance of feature selection for the feature matrix generation.

Main Semantic	Crime flow element code	Description	Suspect 1	Suspect 2
Crime types	HB	House Breaking	0	1
	HK	Hurt by Knife	0	0
	RB	Robbery	1	0
	TH	Theft	0	0
Sub types	ABD/S003	Abduction from the legal guardian	1	0
	ABD/S004	Abducting to marry	0	0
	ABD/S005	Abducting for sexual harassment	0	0
	BGL/S004	Use of stealth	0	1
	BGL/S011	Burglary in business places	0	0
	ROB/S001	Organized vehicle robbery	1	0
Crime flows	BGL/F001	Entering from the window	0	1
	BGL/F002	Entering from the Fanlights	0	0
	BGL/F003	Removing grills	0	1
	BGL/F004	Breaking glasses	0	0
	ROB/F001	Showing identity cards	1	0
	ROB/F003	Wearing uniforms	1	0
	ROB/F004	Robbery using identity cards, uniforms and chains	0	0
	ROB/F009	Seizing inmates	0	0
	ROB/F010	Appearing as CID officers	0	0
	Special category	Retailer 1	Attacking/robbing retailer 1's stores	0
Retailer 2		Attacking/robbing retailer 2's stores	0	0

of the patterns used by the criminals in committing crimes. First, a frequency threshold is generated using characteristic features of the feature matrix at hand which is the matrix of all crimes committed by the same criminal under consideration. The matrix shown in [Table 3](#) is an example to a situation of a feature matrix generated on the previous convictions of a criminal. For the sake of simplicity, let's consider a feature matrix of 10 columns.

If we consider A–J of [Table 3](#) as crime flow features of the corresponding MOs, we can understand that in the first MO the criminal has followed a crime flow of A-E-F-G-I. The same criminal has followed a crime flow of A-D-F-G-I in his second crime. Likewise the other two crime flows are, A-E-F-H-I and A-D-F-G-H-I respectively.

The DMO of a particular criminal is generated using the Apriori method ([Adamo, 2001](#)). Apriori method is used to find the crime entities with the frequency threshold (frt) which is generated according to [Eq. \(2\)](#). A demonstration of the generation of D in [Eq. \(1\)](#) on the properties of feature matrix is shown in [Table 4](#).

$$D = \left\{ d \mid d = \sum_{i=1}^n y_i \right\} \quad (1)$$

$$frt = M_D/n \quad (2)$$



**Table 2** Feature matrix for Suspect 1, generated using the selected modus operandi attributes in Table 1.

	ct1	ct2	ct3	ct4	st1	st2	st3	st4	st5	st6	fl1	fl2	fl3	fl4	fl5	fl6	fl7	fl8	fl9	sc1	sc2
Conviction 1	0	0	1	0	1	0	0	0	0	1	0	0	0	0	1	1	0	0	0	0	0
Conviction 2	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0

**Table 3** Feature matrix generated on four previous convictions of a criminal.

A	B	C	D	E	F	G	H	I	J
1	0	0	0	1	1	1	0	1	0
1	0	0	1	0	1	1	0	1	0
1	0	0	0	1	1	0	1	1	0
1	0	0	1	0	1	1	1	1	0

**Table 4** Column-wise addition of the feature matrix of the suspect under consideration.

A	B	C	D	E	F	G	H	I	J	
1	0	0	0	1	1	1	0	1	0	 Summation 
1	0	0	1	0	1	1	0	1	0	
1	0	0	0	1	1	0	1	1	0	
1	0	0	1	0	1	1	1	1	0	
4	0	0	2	2	4	3	2	4	0	

where,  $D$  = vector of distinct column frequencies of the feature matrix.  $y_i$  = cells in each column,  $M_D$  = Median of  $D$ ,  $n = \sum f$  = number of values or total frequencies,  $c$  = cumulative frequency of the median class,  $h$  = class interval size.

The column-wise addition of the matrix shown in Table 4 gives 4, 0, 0, 2, 2, 4, 3, 2, 4 and 0. The distinct numbers are selected from the resulting vector which results in  $D = [0, 2, 3, 4]$ . The median of  $D$  is then divided by the number of instances (rows) in the matrix as the frt, which is  $2.5/4 = 0.625$  for the above case. Therefore, frt will range from 0 to 1. This value provides an insight to a fair threshold value for the Apriori method to generate the dynamic modus operandi with the most frequent elements. frt is used as the frequency threshold in finding the lengthiest MO with a probability of 0.625 because this value suggests that there is a moderate possibility of one feature having 0.625 probability in each of MO. This results in a dynamic modus operandi (DMO) as shown in Eq. (4), because the only transaction of crime attributes which provides a support of 0.625 is  $\sigma(A, F, G, I)$  as shown in Eq. (3).

$$s = \frac{\sigma(A, F, G, I)}{|T|} = \frac{3}{4} = 0.75 \quad (3)$$

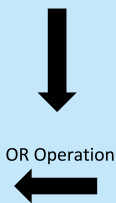
$$\text{DMO} = [1 \ 0 \ 0 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0]. \quad (4)$$

### Generating the complete MO profile (CMOP) of the criminals

The complete MO profile (CMOP) is obtained by the OR operation between the bits of each column of the feature matrix of the corresponding criminal. CMOP guarantees the provision of a composite crime flow by considering all of the previous crime flow entities of a particular criminal. For example, the complete profile for the feature matrix shown in Table 3 is obtained as shown in Table 5.

Therefore,  $\text{CMOP} = [1 \ 0 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0]$ . CMOP contains 1s for each place for which a particular crime flow entity has taken place at least once.

**Table 5** OR operation on the columns to obtain the complete MO profile.

A	B	C	D	E	F	G	H	I	J	
1	0	0	0	1	1	1	0	1	0	
1	0	0	1	0	1	1	0	1	0	
1	0	0	0	1	1	0	1	1	0	
1	0	0	1	0	1	1	1	1	0	
1	0	0	1	1	1	1	1	1	0	

### Finding the deviation probability (DP) of CMOP from the crime MO under consideration (UMO)

First, the deviation of CMOP and UMO is obtained according to Eq. (5). As the binary feature vectors are commonly used to represent patterns, many methods have been invented to find their similarity and distance (Cha, Tappert & Choi, 2010). Euclidean distance, Hamming distance, Manhattan distance, Jaccard, Sorensen and Tanimoto are few of the frequently used measures in that domain (Cha, Tappert & Choi, 2010). This probability value, which is named as the deviation probability (DP), is used to obtain a measurement as to what extent of information is available in the UMO, extra to what is already available in the CMOP of a particular criminal. Let's assume that the bit pattern to be compared with the suspect's modus operandi profile under consideration is  $UMO = [1\ 0\ 0\ 0\ 1\ 1\ 1\ 0\ 0\ 1]$ . Therefore, DP provides the probability of 1s which are available in UMO but not in CMOP.

The deviation probability, DP can be given as,

$$DP = \frac{\sum_{i=1}^n x_i - y_i}{n}, \quad \text{for } x_i = 1, y_i = 0; i = 1, 2, \dots, n \quad (5)$$

where,

$x_i$  = elements of the UMO

$y_i$  = elements of the CMOP.

If we consider the feature matrix on Table 3,

$$\text{Deviation} = [1\ 0\ 0\ 0\ 1\ 1\ 1\ 0\ 0\ 1] - [1\ 0\ 0\ 1\ 1\ 1\ 1\ 1\ 1\ 0]$$

$$\text{Deviation} = [0\ 0\ 0\ -1\ 0\ 0\ 0\ -1\ -1\ 1]. \quad (6)$$

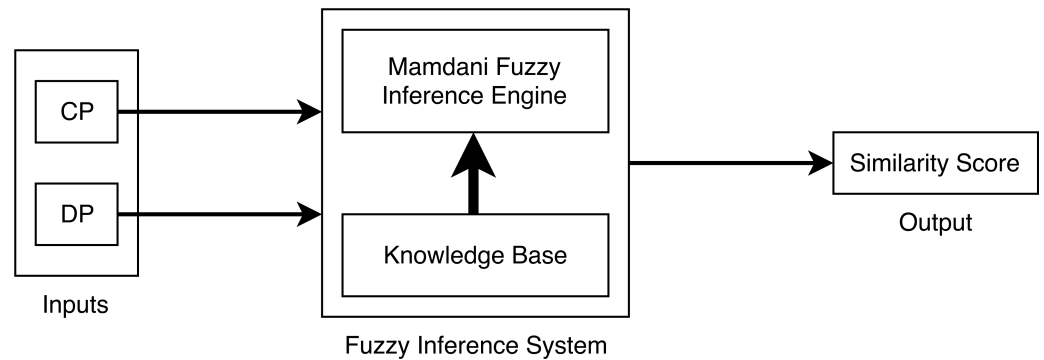
Define  $AD = 1$ , where AD is the number of positive 1s.

Therefore,  $DP = 1/10 = 0.1$ .

As it appears in Expression 6, it produces positive 1s for the places with the features available in UMO but not in CMOP. The higher the DP, higher the amount of extra information available in UMO. Hence, a DP value close to 0 indicates the absence of extra features in UMO.

### Finding the completeness probability (CP) of UMO against DMO

For the same feature matrix which was considered in Table 3, the CP is obtained according to (7). Here, the UMO is compared with DMO to obtain a probability to determine what



**Figure 4** Block diagram of the proposed fuzzy inference system.

extent of features in CP is available in UMO. Therefore, it is derived by the percentage of attributes which are present in both UMO and DMO.

Let  $DMO = \{x_i\}_{i=1}^n$  and  $UMO = \{y_j\}_{j=1}^n$  be two binary sequences.

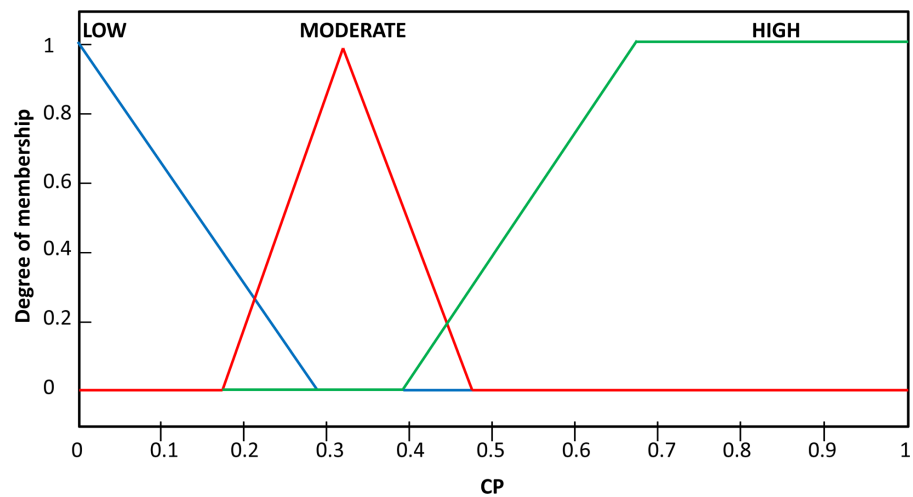
Define,

$$z_k = \begin{cases} 1; & x_i = y_j \\ 0; & \text{otherwise} \end{cases} \quad \text{Then, } CP = \frac{\sum_{k=1}^n z_k}{n} \text{ is the completeness probability.} \quad (7)$$

For example, if we consider  $DMO = [1\ 0\ 0\ 0\ 0\ 1\ 1\ 0\ 1\ 0]$ , then for the  $UMO = [1\ 0\ 0\ 0\ 1\ 1\ 1\ 0\ 0\ 1]$  a CP of  $3/10 = 0.3$  is generated as in the 1st, 6th and 7th positions there are ones in both DMO and UMO. The higher the CP value, the more the UMO is composed of crime flow entities which are available in the DMO. Therefore, a CP value close to 1 indicates that the completeness of UMO compared to DMO is 100%.

### Building a fuzzy inference system to obtain the final similarity score

The vagueness of the two measurements CP and DP generates a difficulty in calculating a similarity score using crisp logic. Therefore, the two parameters CP and DP were adapted into a fuzzy inference system which accepts two inputs and provides a score for the similarity between a suspect and a crime. Figure 4 shows a block diagram of the proposed fuzzy inference system. Mamdani fuzzy inference was used as an attempt to solve a control problem by a set of linguistic rules obtained from experienced human operators (Mamdani & Assilina, 1975). First, the rule base of the fuzzy controller was defined by observing the variations of CP and DP. The membership functions of the inputs and outputs were then adjusted in such a way that the parameters which seem to be wrong can be fine-tuned, which is a common practice in defining fuzzy inference systems (Godjevac, 1997). Literature shows many methods used in fine tuning the fuzzy parameters. Usage of adaptive networks (Sun, 1994) and Neuro-fuzzy systems (Abraham, Nath & Mahanti, 2001) in fine tuning the fuzzy parameters have received more attention. The problem at our hand was to generate a fuzzy inference system which generates the highest similarity score when the DP value goes down and CP value goes up. We conducted a manual mapping procedure for the fuzzy membership functions. Therefore, the input and output space of the two inputs CP and DP and the output were partitioned into 3 subsets. Namely,



**Figure 5** Input fuzzy variable 1: CP.

LOW, MODERATE and HIGH. Center of gravity was used as the defuzzification strategy of the fuzzy controller. Mamdani fuzzy inference was especially selected for the similarity score generation procedure, for the highly intuitive knowledge base it offers due to the fact that both antecedents and the consequents of the rules are expressed as linguistic constraints (Zadeh, 1997). First, we selected all of these membership functions with 50% overlap. Then the tuning procedure was conducted during which we adjusted either the left and/or right spread and/or overlapping to get the best possible similarity score for the given DP and CP. This procedure was conducted until the FIS generated satisfactory results.

Figures 5 and 6 show the fuzzy inputs of the Fuzzy Inference System (FIS) which correspond to CP and DP values respectively. Figure 7 depicts the fuzzy output of the FIS. As the Figs. 5–7 depict, all the different levels of membership functions under each input and the output are selected to be triangular and trapezoidal functions as triangular or trapezoidal shapes are simple to implement and computationally efficient (MathWorks Inc, 1994–2015). As shown in Fig. 7, the universe of discourse of similarity score (fuzzy output) ranges from 0 to 100. The defuzzified score which is generated from the FIS is considered as the measurement for how close the modus operandi under consideration is to a particular suspect's profile. A higher score value close to 100 provides a good indication about a high similarity between the modus operandi of the crime and suspect under consideration.

The fuzzy rule derivation of the fuzzy controller is heuristic in nature. According to the calculations of the two inputs, higher values of CP, close to 1 and lower values of DP close to 0, positively affect the final similarity score. The rule base of the fuzzy model is generated accordingly. The rule base provides a non-sparse rule composition of 9 combinations as illustrated in Fig. 8.

The rule surface of the fuzzy controller depicted in Fig. 9, shows the variation of the similarity score with the changes of the two inputs CP and DP. According to the figure it's perfectly visible, for higher values of CP (close to 1) and for lower values of DP (close to 0), the fuzzy controller generates higher values for the similarity score which are close to 100.

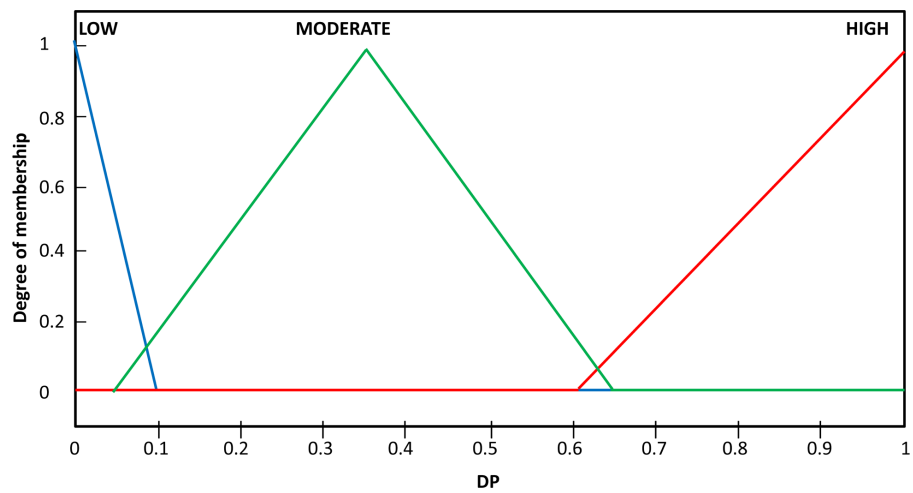


Figure 6 Input fuzzy variable 2: DP.

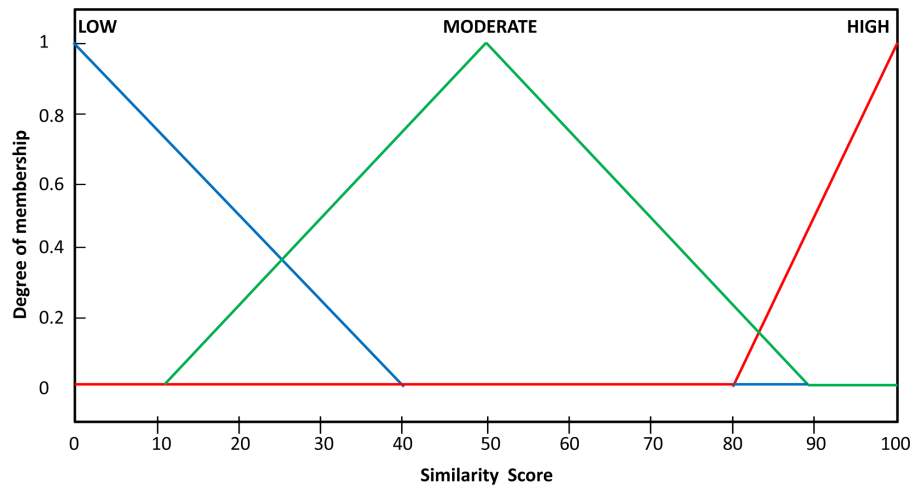
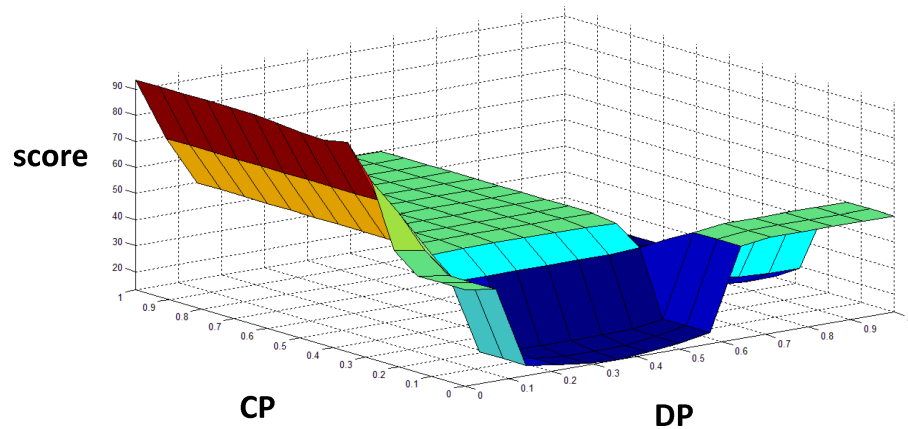


Figure 7 Output fuzzy variable: similarity score.

CP \ DP	LOW	MODERATE	HIGH
LOW	MODERATE	LOW	MODERATE
MODERATE	HIGH	MODERATE	LOW
HIGH	HIGH	MODERATE	LOW

Figure 8 Fuzzy rule set of the rule base of the inference system.



**Figure 9** Rule surface of the fuzzy controller.

### Classification of the UMO under the class with the highest similarity score

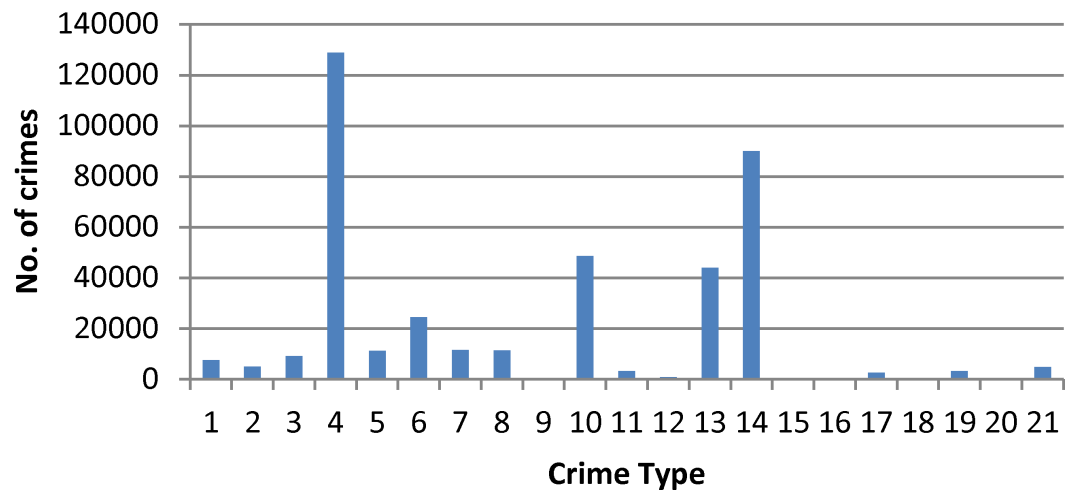
When the algorithm is used to find associations between modi operandi of criminals and modi operandi of crimes, the similarity score which is generated from the newly proposed method can be used directly. A similarity score which is close to 100 would suggest that the criminal has a very high tendency to have committed the crime which is under investigation. Therefore, the similarity scores can be used to classify a particular modus operandi to a most probable suspect with the highest similarity score.

The proposed method was developed by using MATLAB 7.12.0 (R2011a) (*MathWorks, 1994–2015a*). All the necessary implementations were conducted using the MATLAB Script editor (*MathWorks, 1994–2015b*) apart from the FIS which was implemented using the MATLAB fuzzy toolbox (*MathWorks, 1994–2015c*). The nine classification algorithms which were used for the performance comparison were classification algorithms which are already packaged with the WEKA 3.6.12 tool (*Hall et al., 2009*).

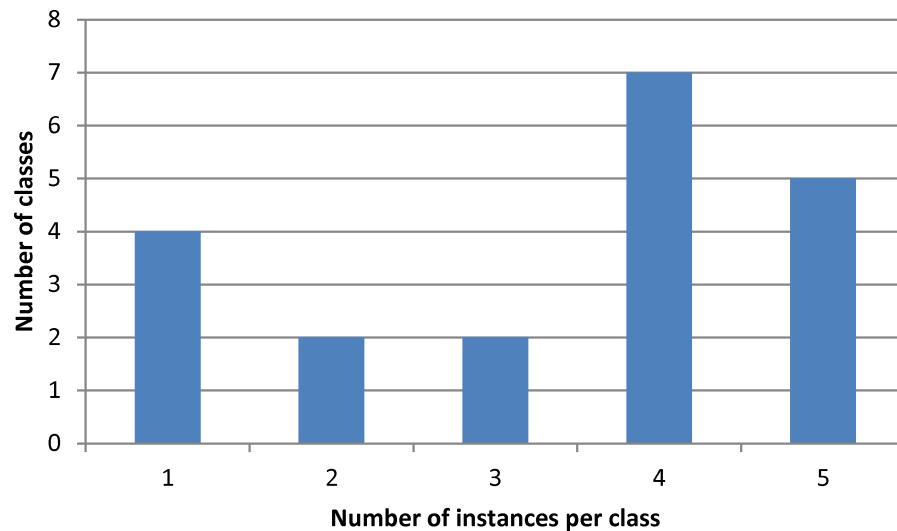
## RESULTS AND DISCUSSION

The method was tested with a crime data set obtained from Sri Lanka Police. [Figure 10](#) shows the crime frequencies in Sri Lanka by the crime types from 2005 to 2011. It shows only 21 crime types because the five new crime types were introduced in 2015. 4th column denoting House Breaking and Theft shows the highest number of occurrences. 14: Theft of property, 10: Robbery, 13: Cheating/Misappropriation, 6: Hurt by Knife, 7: Homicide, 8: Rape/Incest, 5: Grievous Hurt, 3: Mischief over Rs. 5,000/=, 1: Abduction/Kidnapping comes next. For the validation of the algorithm, 7 crime types out of these 10 types were selected for the testing data set. They are, House Breaking and Theft, Theft of property, Robbery, Homicide, Rape/Incest, Grievous Hurt, Abduction/Kidnapping. A total of 31 crime flows were selected which are common to the seven selected crime types. The data set is also composed of eight sub types and two special categories. Altogether the data set consisted of 67 instances in which each instance is composed of 48 attribute values. The data set is distributed over 20 classes (criminals).





**Figure 10** Frequency of different crime types from year 2005–2011.



**Figure 11** Distribution of modus operandi instances over the classes of the dataset.

All the tests were performed in a Windows computer with Intel (R) Core (IM) i7-2670QM CPU of 2.20 GHz and a RAM of 8 GB. The histogram of the instance distribution over the classes is shown in Fig. 11.

A 10 fold cross validation (Refaeilzadeh, Tang & Liu, 2009) was used on the data set for a fair testing procedure. In 10-fold cross validation, the data set is divided into 10 subsets, and the holdout method is repeated 10 times. Each time, one of the 10 subsets is used as the test set and the other nine subsets are put together to form a training set. Then the average error across all 10 trials is computed (Refaeilzadeh, Tang & Liu, 2009).

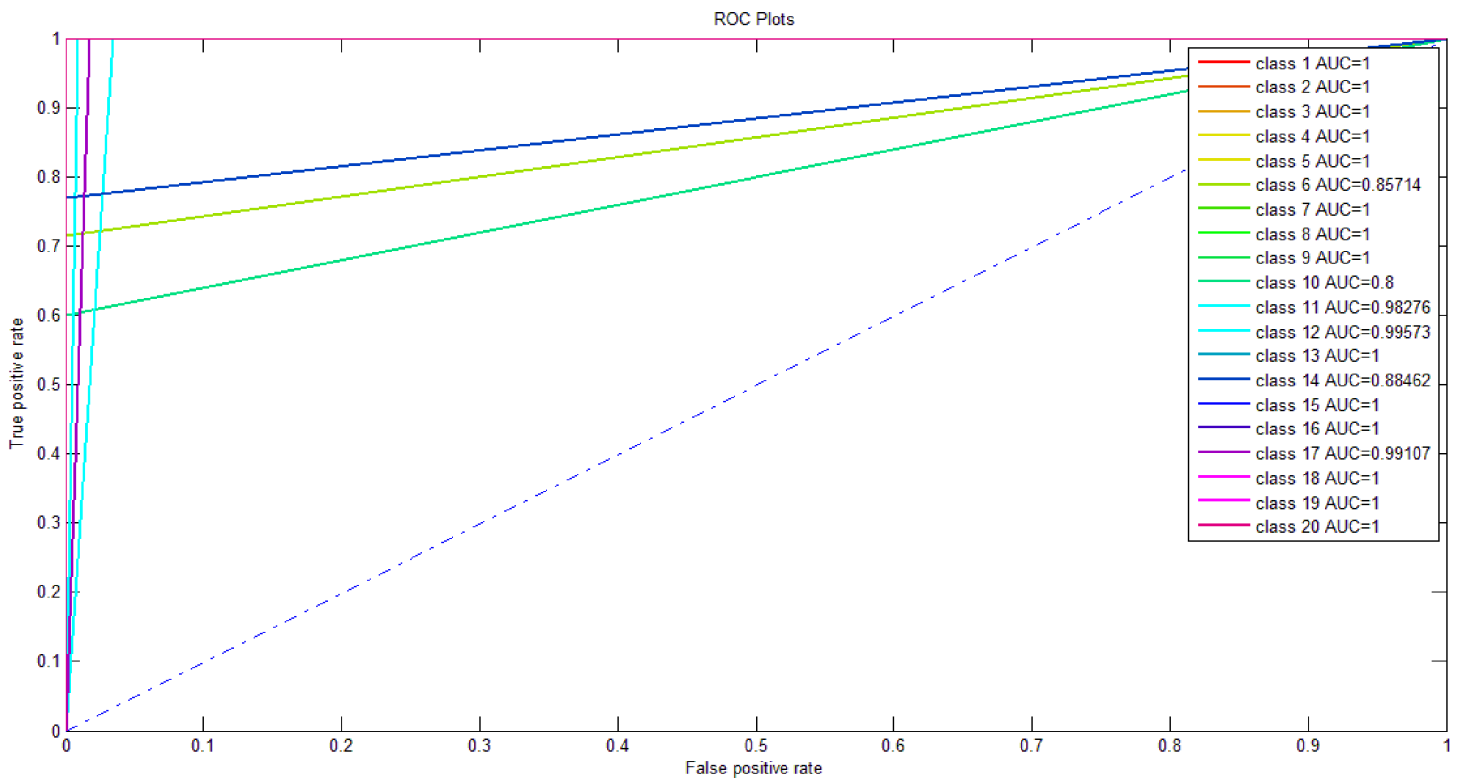
The test results of modus operandi classifications in Area Under Curve (AUC) (Hanley & McNeil, 1982), and time elapsed for the classification are shown in Table 6. A Receiver Operating Characteristic (ROC) curve is a two dimensional graphical illustration of the trade-off between the true positive rate (sensitivity) and false positive rate (1-specificity).

**Table 6** Results returned by the fuzzy based binary feature profiling for the modus operandi analysis on actual data.

Data set (number)	Oversampling or Under-sampling value	AUC	Average time elapsed
1	2	0.5417	0.0015
2	3	0.5562	0.0011
3	4	0.5965	0.0014
4	N/A	0.6937	0.0010
5	5	0.6612	0.0011
6	6	0.7063	0.0011
7	10	0.8033	0.0012
8	20	0.9339	0.0013
9	30	0.9661	0.0014
10	40	0.9637	0.0015
11	50	0.9756	0.0016
12	60	0.9626	0.0018
13	70	0.9365	0.0019
14	80	0.9391	0.0023
15	90	0.9671	0.0029

Figure 12 depicts the ROC curve plotted on the classification results obtained by the newly proposed method on the crime data set. In the particular instance which is shown in Fig. 12, all the ROC curves related to the crime data set are plotted well over the diagonal line and all of them have returned AUC values which are either equal to 1 or very close to 1, providing a very good classification.

To prepare the data set which was used under this research, a crime data set of around 3,000 instances was analyzed. Due to the limitations of the real crime data set, it was quite a complex task to prepare a data set with a collection of sufficient modus operandi where each instance has a considerable flow of crime flows. Therefore, only a sample of 67 instances could be filtered from the population to generate a representative data set and it was verified by a domain expert before being used in the analysis. As the number of instances was around 67, it can be considered as an under-represented data sample. Another reason for the data set to become under-represented was the challenge in finding classes/criminals with more than one crime committed. The actual crime data set which is used for the testing purposes is imbalanced as it is apparent in Fig. 11. This imbalanced nature of the data set may produce biased results. To make the classification process unbiased, we used the concept of oversampling. Oversampling and under-sampling are two concepts which are used in overcoming class imbalance problems in input data sets. Oversampling and under-sampling are two different categories of resampling approaches, where in oversampling the small classes are incorporated with repeated instances to make them reach a size close to larger classes, whereas in under-sampling, the number of instances is decreased in such a way that the number of instances reach a size close to the smaller classes (Estabrooks, Jo & Japkowicz, 2004).



**Figure 12** ROC curves returned by the newly proposed method on the 20 classes of crime data set.

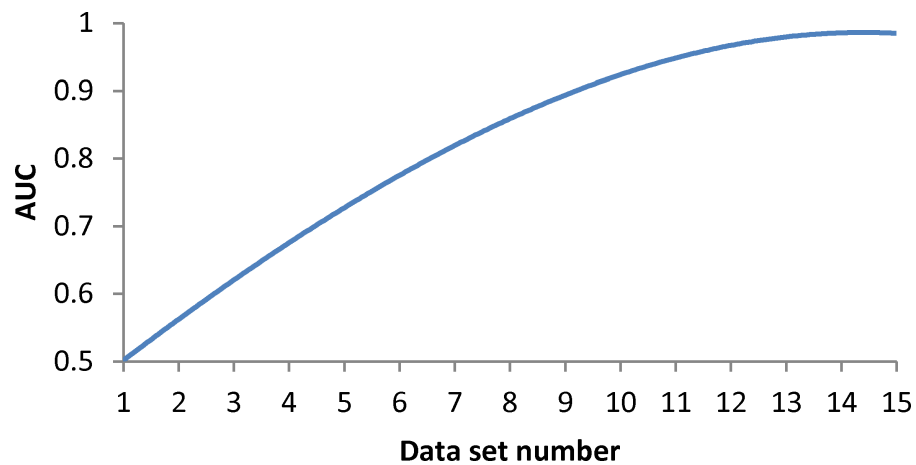
Table 6 shows the results returned by the fuzzy based binary feature profiling which was conducted on the actual crime data set. As shown in the table, there is an increase in the accuracy when the input data set undergoes oversampling. Since the maximum number of instances available under one suspect is equal to 5, under-sampling does not provide a good accuracy. The results prove that the new algorithm works well for a balanced data set as the new method showed an increase in performance when the data set is subjected to an oversampling greater than or equal to 5.

Figure 13 shows the change in x AUC with the increase of sampling which starts from under-sampling of 2 and goes on to an over sampling of 90. According to the plot it can be observed that the AUC values are increased when the oversampling is increased.

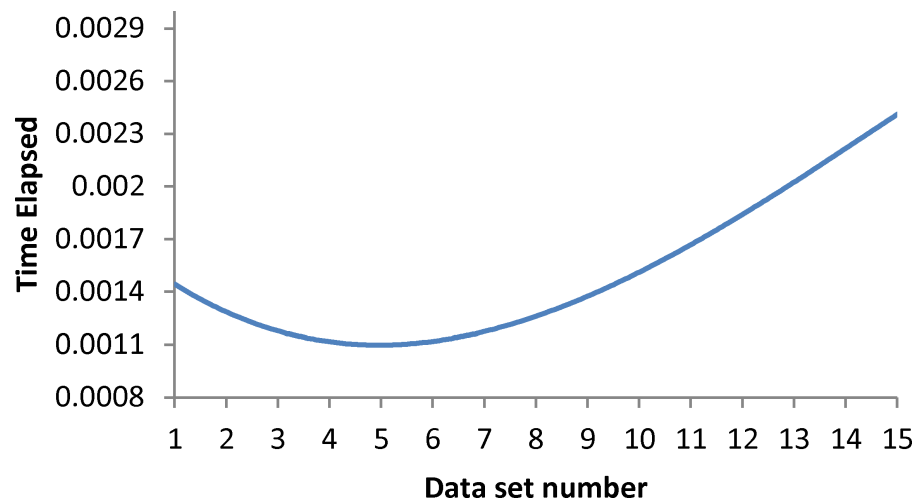
The execution time of the algorithm was 0.001 s when there is no oversampling or under-sampling. The maximum execution time is 0.0031 when there is an oversampling of 90. According to the plot shown in Fig. 14, it is clear that there is an increase of execution time as the oversampling size increases. But, the overall execution time is always remained under 3 ms.

### Overview of the classification algorithms used for the comparison

It is a known fact that there is no single algorithm which can be categorized as the best to solve any problem. Different classification algorithms may perform differently in different situations (Wolpert, 1996). Therefore, the newly proposed method was tested against ten other open classification data sets (The information about these data sets is provided in



**Figure 13** Change of AUC values with oversampling.



**Figure 14** Change of time elapsed for the 15 data sets.

Table 7) and the performance was evaluated against the results obtained from nine other well-known classification techniques, thereby assessing the quality of the newly proposed method. The nine other classification algorithms include, Logistic Regression, J48 Decision Tree, Radial Basis Function Network (RBFNetwork), Multi-Layer Perceptron (MLP), Naive Bayes Classifier, Sequential Minimal Optimization (SMO) algorithm, KStar instance based classifier, Best-first decision tree (BFTree) classifier, and Logistic Model Tree (LMT) classifier. These classifiers represent four classes of classification algorithms. Namely, function based classifiers, Tree based classifiers, Bayesian classifiers and Lazy classifiers.

Logistic Regression learns conditional probability distribution. Relating qualitative variables to other variables through a logistic cumulative distribution functional form is logistic regression (*Chang & Lin, 2004*). J48 is an open source java implementation of the C4.5 decision tree algorithm (*Machine Learning Group at the University of Waikato*). A decision tree consists of internal nodes that specify tests on individual input variables

**Table 7** Description of the classification data sets for performance comparison.

Data set	Description	Number of instances	No of attributes
Dermatology data set ( <i>Ilter &amp; Guvenir, 1998</i> )	This database has been created on a dermatology test carried out on skin samples which have been taken for the evaluation of 22 histopathological features. The values of the histopathological features have been determined by an analysis of the samples under a microscope. In the dataset constructed for this domain, the family history feature has the value 1 if any of these diseases has been observed in the family, and 0 otherwise. Every other feature (clinical and histopathological) was given a degree in the range of 0 to 3. Here, 0 indicates that the feature was not present, 3 indicates the largest amount possible, and 1, 2 indicate the relative intermediate values.	336	33
Balance scale data set ( <i>Hume, 1994</i> )	This data set has been generated to model psychological experimental results. Each example is classified as having the balance scale tip to the right, tip to the left, or be balanced. The attributes are the left weight, the left distance, the right weight, and the right distance. The correct way to find the class is the greater of (left-distance * left-weight) and (right-distance * right-weight). If they are equal, it is balanced. There are 3 classes (L, B, R), five levels of Left-Weight (1, 2, 3, 4, 5), five levels of Left-Distance (1, 2, 3, 4, 5), five levels of Right-weight (1, 2, 3, 4, 5) and five levels of Right-Distance (1, 2, 3, 4, 5).	625	4
Balloons data set ( <i>Pazzani, 1991a</i> )	This data set has been generated using an experiment of stretching a collection of balloons carried out on a group of adults and children ( <i>Pazzani, 1991b</i> ). In the data set, Inflated is true if (color = yellow and size = small) or (age = adult and act = stretch). In the data set there are two main output classes, namely T if inflated and F if not inflated, two colors yellow and purple, two sizes, large and small, two act types, stretch and dip, and two age groups, adult and child.	20	4
Car evaluation data set ( <i>Bohanec &amp; Zupan, 1997a</i> )	Car Evaluation Database has been derived from a simple hierarchical decision model originally developed for the demonstration of DEX by <i>Bohanec &amp; Rajkovic (1990)</i> . The Car Evaluation Database contains examples with information that is directly related to CAR. They are buying, maint, doors, persons, lug_boot and safety. The attribute buying is the buying price which is considered to have four levels v-high, high, med, low. Maint is the price of the maintenance which contains the four levels, v-high, high, med, low. Doors have the four levels 2, 3, 4, 5-more. Person (capacity in terms of persons to carry), lug_boot (the size of luggage boot) and safety (estimated safety of the car) have 3 levels each.	1,728	6
Soybean data set ( <i>Fisher, 1987; Michalski, 1980</i> )	This is a small subset of the original soybean database. The data set is distributed over four classes, D1, D2, D3 and D4. The 35 categorical variables represent different levels of qualities of the soybean vegetable. These categorical variables include, plant-stand, precip, temp, hail, crop-hist, area-damaged, severity, seed-tmt, germination, lant-growth, leaves, leafspots-halo, leafspots-marg, leafspot-size, leaf-shread, leaf-malf, leaf-mild, stem, lodging, stem-cankers, canker-lesion, fruiting-bodies, external, mycelium, int-discolor, sclerotia, fruit-pods, fruit, seed, mold-growth, seed-discolor, seed-size, shriveling and roots. The number of levels represented by each variable varies from 2 to 3.	47	35
Lenses data set ( <i>Julien, 1990</i> )	Lenses data set is a small database about fitting contact lenses. The data set is composed of five attributes including the class variable. The data set has three classes. Age of the patient, spectacle prescription, astigmatic, tear production rate are the attributes of the data set. The attributes contain at least of two categories and at most of three categories.	24	4

(continued on next page)

Table 7 (continued)

Data set	Description	Number of instances	No of attributes
Nursery data set (Bohanec & Zupan, 1997b)	Nursery Database has been derived from a hierarchical decision model originally developed to rank applications for nursery schools. It has been used during several years in 1980s when there has been excessive enrollment to these schools in Ljubljana, Slovenia, and the rejected applications frequently needed an objective explanation. The final decision depended on three sub problems: occupation of parents and child's nursery, family structure and financial standing, and social and health picture of the family. The model has been developed within expert system shell for decision making (Bohanec & Rajkovic, 1990).	12,960	8
Tic-tac-toe data set (Aha, 1991)	This database encodes the complete set of possible board configurations at the end of tic-tac-toe games, where "x" is assumed to have played first. The target concept is "win for x" (i.e., true when "x" has one of 8 possible ways to create a "three-in-a-row").	958	9
SPECT heart data set (Kurgan & Cios, 2001)	The data set describes diagnosis of cardiac Single Proton Emission Computed Tomography (SPECT) images. Each of the patients is classified into two categories: normal and abnormal. The database of 267 SPECT image sets (patients) were processed to extract features that summarize the original SPECT images. The instances are described by 23 binary attributes including the class variable.	267	22
MONK's problems data set (Thrun, 1992)	The MONK's problems have been the basis of a first international comparison of learning algorithms. The result of this comparison is summarized in "The MONK's Problems" (Thrun et al., 1991). There are three MONK's problems. The domains for all MONK's problems are the same. The data set is composed of 7 attributes and a binary class variable.	432	7

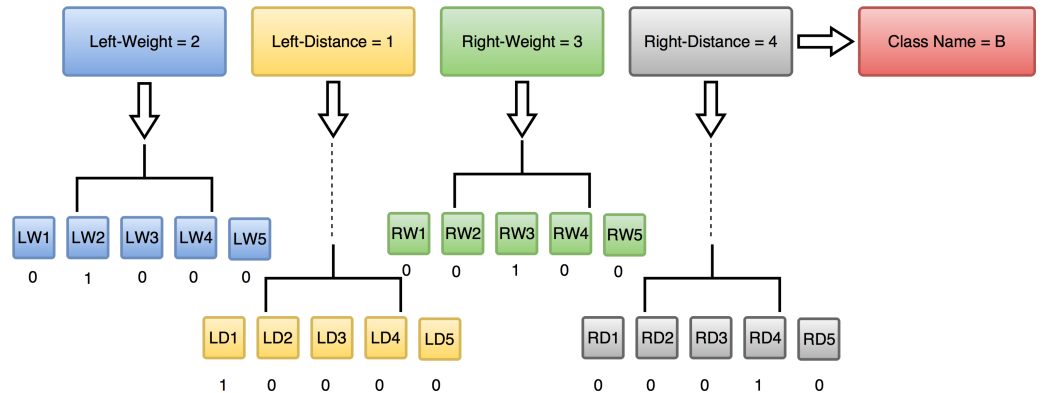
or attributes that split the data into smaller subsets, and a series of leaf nodes assigning a class to each of the observations in the resulting segments. The C4.5 algorithm constructs decision trees using the concept of information entropy (Quinlan, 1993). Neural networks are flexible in being modeled virtually for any non-linear association between input variables and target variables (Bishop, 1995). Both Radial basis function networks and MLP networks are neural networks (Jayawardena, Fernando & Zhou, 1997). Bayesian classifiers assign the most likely class to a given example described by its feature vector (Rish, 2001). SMO is an implementation of John Platt's sequential minimal optimization algorithm for training a support vector classifier. It globally replaces all missing values and transforms nominal attributes into binary one. It also normalizes all attributes by default (Platt, 1999; Keerthi et al., 2001). KStar (K\*) is an instance-based classifier which uses an entropy-based distance function (Cleary & Leonard, 1995). BFTree uses binary split for both nominal and numeric attributes (Friedman, Hastie & Tibshirani, 2000). LMT is a classifier for building 'logistic model trees', which are classification trees with logistic regression functions at the leaves (Landwehr, Hall & Frank, 2005; Sumner, Frank & Hall, 2005).

As the newly proposed method accepts only binary input variables, the data sets which are used for the analysis must be preprocessed into the acceptable format. For example, the "balance scale" data set is composed of four attributes. Table 8 shows the attributes and their information of the balance scale data set.

Therefore, the data set was adjusted as shown in Fig. 15, prior to using it with the proposed method. Each category of a particular attribute is represented by a dummy variable. For example, Left-Weight attribute results in five attributes in the preprocessed

**Table 8** Attribute information of the balance data set.

Attribute	Number of categories	Categories
Class name	3	L, B, R
Left-weight	5	1, 2, 3, 4, 5
Left-distance	5	1, 2, 3, 4, 5
Right-weight	5	1, 2, 3, 4, 5
Right-distance	5	1, 2, 3, 4, 5

**Figure 15** Schematic diagram used for pre-processing of the balance dataset in such a way that it matches the format of inputs of the newly proposed method.

data set and each attribute is represented using five binary variables as LW1, LW2, LW3, LW4 and LW5 where the presence of the attribute denotes 1 and 0 otherwise. As depicted in Fig. 15, if Left-Weight has a value of 2 in an instance it results in 1 for the corresponding derived attribute that is LW2. Therefore, if there is an instance where Left-Weight = 2, Left-Distance = 1, Right-Weight = 3 and Right-Distance = 4, Class Name = B, it is represented as LW1 = 0, LW2 = 1, LW3 = 0, LW4 = 0, LW5 = 0, LD1 = 1, LD2 = 0, LD3 = 0, LD4 = 0, LD5 = 0, RW1 = 0, RW2 = 0, RW3 = 1, RW4 = 0, RW5 = 0, RD1 = 0, RD2 = 0, RD3 = 0, RD4 = 1, RD5 = 0, Class Name = B.

The pre-processed data is then fed to the newly proposed algorithm and the nine other algorithms. Performances were compared based on AUC analysis of the ROC curves, and the processing time for model generation. 10 fold cross validation was used under each test for fair testing procedure. For the sake of simplicity, the newly proposed modes operandi analysis algorithm was acronymed as BFPM (Binary feature profiling methodology).

As all the data sets which were used for the tests are composed of multi classes, weighted average AUC was used, where each target class is weighted according to its prevalence as given in Eq. (8). Weighted average was used in order to prevent target classes with smaller instance counts from adversely affecting the results (Hempstalk & Frank, 2008).

$$AUC_{weighted} = \sum_{\forall c_i \in C} AUC(c_i) \times p(c_i) \quad (8)$$

Table 9 shows the weighted average AUC values obtained for each data set under each classification algorithm.



**Table 9** Weighted average AUC values obtained by the algorithms on classifying the data sets.

	<b>BFPM</b>	<b>Logistic regression</b>	<b>J48</b>	<b>Radial basis function network</b>	<b>Multi-layer perceptron</b>	<b>Naive Bayes classifier</b>	<b>SMO</b>	<b>KStar</b>	<b>BFTree</b>	<b>LMT</b>
Dermatology data set	1	0.9990	0.9750	0.9860	0.9980	0.9980	0.9930	0.9970	0.9690	0.9960
Balance scale data set	0.7945	0.9760	0.8110	0.9680	0.9770	0.9710	0.8830	0.9510	0.8130	0.9810
Balloons data set	1	1	1	1	1	1	1	1	1	1
Car evaluation data set	0.8087	0.9900	0.9760	0.9740	1	0.9760	0.9550	0.9970	0.9940	0.9990
Soybean data set	1	1	0.9860	1	1	0.9760	1	1	0.9740	1
Lenses data set	0.9537	0.7470	0.8400	0.9170	0.8390	0.8700	0.7250	0.8870	0.8670	0.7980
Nursery data set	0.9100	0.9880	0.9950	0.9870	1	0.9820	0.9640	0.9980	0.9990	0.9990
Tic-tac-toe data set	0.9167	0.9960	0.8970	0.7340	0.9940	0.7440	0.9760	0.9990	0.9450	0.9920
SPECT heart data set	0.7857	0.8310	0.7560	0.8400	0.7860	0.8490	0.7070	0.7850	0.7230	0.8410
MONK's problems data set	0.8333	0.7050	0.9940	0.8130	0.9980	0.7120	0.7460	0.9970	0.9550	0.9880

**Table 10** Friedman’s mean rank values returned on the data available in Table 9.

Method	Mean rank
MLP	7.70
LMT	7.10
KStar	6.95
LogisticRegression	5.95
RBFNetworks	5.05
NaiveBayesClassifier	5.05
<b>BFBM</b>	<b>4.95</b>
BFTree	4.50
J48	4.20
SMO	3.55

Friedman’s rank test is a nonparametric test analogous to a standard one-way repeated-measures analysis of variance (Howell, 2013). The Friedman’s rank test results returned on the AUC test data are shown in Table 10. This test returns a test statistic ( $\chi^2$ ) value (“Chi-square”) of 21.339, degree of freedom of 9 and a  $p$ -value of 0.011, proving that there is an overall statistically significant difference between the mean ranks of the classification algorithms. According to the table, the highest mean rank is returned for MLP while the lowest mean rank is returned for SMO, proving that MLP provides the best performance while SMO provides the least performance for the 10 data sets tested. Therefore, it indicates that the new model provides a better performance than BFTree, J48 and SMO algorithms for the 10 data sets tested.

The average processing times elapsed for each algorithm to classify the data sets are given in Table 11. Friedman’s rank test on the data of Table 11 returned the results shown in Table 12 in which the mean rank values prove better efficiency for the new method than J48, LogisticRegression, SMO, RBFNetworks, BFTree, MLP and LMT. The test statistic ( $\chi^2$ ) value (“Chi-square”) of 73.058, degree of freedom of 9 and a  $p$ -value of 0.000, proves that there is an overall statistically significant difference between the mean ranks of the classification algorithms.

Friedman’s rank test results for the two measurements, AUC and time elapsed conclude that the newly proposed method provides acceptable results against the nine other well established classification algorithms.

## CONCLUSION

The studies of modus operandi help crime investigation by letting the police officers to solve crimes by linking suspects to crimes. Though there are many descriptive studies available under modus operandi analysis, only a small amount of work is available under computer science. Many of these methods have been derived using the methods based on link analysis. However, the accuracy of these methods is always compromised due to the cognitive biases of the criminals.

A novel Fuzzy based Binary Feature Profiling method (BFBM) to find associations between crimes and criminals, using modus operandi is introduced. The newly proposed

**Table 11** Average processing time for each algorithm on the classification of the ten data sets.

	BFPM	Logistic regression	J48	Radial basis function network	Multi-layer perceptron	Naive Bayes classifier	SMO	KStar	BFTree	LMT
Dermatology data set	0.0027	0.3900	0.0800	0.3800	2.7300	0.0500	0.1400	0.2800	0.1100	2.9000
Balance scale data set	0.0030	0.0300	0.0800	0.2200	0.4800	0	0.0500	0	0.0600	0.8400
Balloons data set	0	0	0	0	0.0200	0	0.0200	0	0.0200	0.0500
Car evaluation data set	0.0048	0.4200	0.0300	0.2700	13.4000	0.0200	0.1900	0	0.4800	13.3400
Soybean data set	0.0042	0.0500	0.0700	0.1800	0.4300	0	0.0500	0	0.2300	0.9200
Lenses data set	0.0009	0	0	0.0300	0.0800	0	0.0300	0	0.0100	0.0300
Nursery data set	0.0013	8.8600	0.2500	16.4300	127.8600	0.0300	23.0300	0	7.9900	240.4600
Tic-tac-toe data set	0.0091	0.1900	0.0100	0.1300	18.1800	0.0100	0.6000	0	0.6100	54.6700
SPECT heart data set	0.0073	0.0600	0.0100	0.0700	3.8800	0	0.0400	0	0.2300	2.1500
MONK's problem data set	0.0058	0.0800	0.0100	0.0600	5.5100	0	0.1300	0	0.2100	1.7200

**Table 12** Mean rank values returned by the Friedman's rank test on the time values available in Table 11.

Method	Mean rank
KStar	2.10
NaiveBayesClassifier	2.35
<b>BFBM</b>	<b>2.75</b>
J48	4.15
LogisticRegression	5.35
SMO	6.25
RBFnetworks	6.35
BFTree	6.90
MLP	9.30
LMT	9.50

method subjects not only the properties of the present, but also the properties of his/her previous convictions. The concept of dynamic modus operandi which is available in the proposed method considers the modi operandi of all of his/her previous convictions to provide a fair rectification to the errors which result due to the human cognition. Dynamic MO uses frequent item set mining to result in a generalized binary feature vector. Complete MO profile also encapsulates past modi operandi of a particular criminal by aggregating the modi operandi of all of his/her previous convictions to one binary feature vector. This feature also guarantees a usage of criminal's past crime record with more generalizability. Completeness probability measures how much information is available in the new crime which is not available in the complete MO profile. Therefore, this measurement provides the capability of measuring how much extra amount of information is carried by the MO of the new crime. The deviation probability provides a notion about how much the new MO deviates from the most frequent crime flows which are available in the dynamic MO of a particular criminal. The vagueness and the impreciseness prompted the fact that it is not possible to use crisp logic to generate the similarity score. Therefore, a fuzzy inference system was modeled to generate the similarity score.

Due to the under-represented and imbalanced properties of the actual data set, the new method has returned a lower performance when it was proposed to the data set without any rectification on the data set. However, with the introduction of over sampling, the method returned a very good performance, allowing one to arrive at the conclusion that the method could provide acceptable results for a balanced data set. The method generated favorable results in providing a good similarity measurement to suggest the connections between crimes and criminals. The fuzzy controller of the new approach guarantees to resemble the human reasoning process by confirming the usage of human operator knowledge to deal with nonlinearity of the actual situation. The newly proposed method was then adapted into a classification algorithm for the validation and comparison with other classification algorithms. The comparison of the new method with the well-established classification algorithms confirmed the generalizability of the new method.

The method only provides the capability to process the categorical data sets. If there are any continuous variables in the data set, the values must be introduced with categories

before further processing. The method can be further extended to directly accept the continuous attributes. As the center of gravity method is used for the defuzzification process, further optimizations can be done by simplifying the defuzzification procedure. Adapting the fuzzy inference engine to a Sugeno (*Takagi & Sugeno, 1985*) type and converting the defuzzification method to a more computationally efficient method such as the weighted average (*Wu & Mendel, 2007*) method would provide a less complex computation. This would result in even less processing time when the sophistication of the data set rises.

## ADDITIONAL INFORMATION AND DECLARATIONS

### Funding

This work was funded by the National Research Council (NRC) of Sri Lanka (Grant number: 11-071). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

### Grant Disclosures

The following grant information was disclosed by the authors:  
National Research Council (NRC): 11-071.

### Competing Interests

The authors declare there are no competing interests.

### Author Contributions

- Mahawaga Arachchige Pathum Chamikara conceived and designed the experiments, performed the experiments, analyzed the data, contributed reagents/materials/analysis tools, wrote the paper, prepared figures and/or tables, performed the computation work.
- Akalanka Galappaththi conceived and designed the experiments, contributed reagents/materials/analysis tools.
- Roshan Dharshana Yapa and Ruwan Dharshana Nawarathna conceived and designed the experiments, performed the experiments, performed the computation work, reviewed drafts of the paper.
- Saluka Ranasinghe Kodituwakku conceived and designed the experiments, contributed reagents/materials/analysis tools, reviewed drafts of the paper.
- Jagath Gunatilake conceived and designed the experiments, reviewed drafts of the paper.
- Aththanapola Arachchilage Chathranee Anumitha Jayathilake conceived and designed the experiments, performed the experiments, analyzed the data, contributed reagents/materials/analysis tools, performed the computation work.
- Liwan Liyanage conceived and designed the experiments, analyzed the data, contributed reagents/materials/analysis tools, reviewed drafts of the paper.

### Data Availability

The following information was supplied regarding data availability:

The data sets were taken from: <https://archive.ics.uci.edu/ml/datasets.html>. The individual links to separate data sets are provided inline with the text of the paper.

## Supplemental Information

Supplemental information for this article can be found online at <http://dx.doi.org/10.7717/peerj-cs.65#supplemental-information>.

## REFERENCES

- Abraham A, Nath B, Mahanti PK. 2001.** Hybrid intelligent systems for stock market analysis. In: *Computational science-ICCS. Lecture notes in computer science*, vol. 2074. 337–345.
- Adamo JM. 2001.** *Data mining for association rules and sequential patterns, Sequential and Parallel Algorithms*. 1st edition. New York: Springer Science & Business Media.
- Agrawal R, Imielinski J, Swami A. 1993.** mining association rule between sets of items in large databases. In: *Proceedings of the ACM SIGMOD international conference of management of data*. New York: ACM, 207–216.
- Aha DW. 1991.** UCI Machine Learning Repository. Available at <https://archive.ics.uci.edu/ml/datasets/Tic-Tac-Toe+Endgame>.
- Bennell C, Canter DV. 2002.** Linking commercial burglaries by modus operandi: tests using regression and ROC analysis. *Science & Justice* **42(3)**:153–164 DOI 10.1016/S1355-0306(02)71820-0.
- Bennell C, Jones NJ. 2005.** Between a ROC and a hard place: a method for linking serial burglaries by modus operandi. *Journal of Investigative Psychology and Offender Profiling* **2(1)**:23–41 DOI 10.1002/jip.21.
- Berry MJ, Linoff G. 2011.** *Data mining techniques: for marketing, sales, and customer support*. 3rd edition. Hoboken: Wiley.
- Bishop CM. 1995.** *Neural networks for pattern recognition*. Oxford: Oxford University Press.
- Bohanec M, Rajkovic V. 1990.** Expert system for decision making. *Sistemica* **1(1)**:145–157.
- Bohanec M, Zupan B. 1997a.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Car+Evaluation>.
- Bohanec M, Zupan B. 1997b.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Nursery>.
- Borg A, Boldt M, Lavesson N, Melander U, Boeva V. 2014.** Detecting serial residential burglaries using clustering. *Expert Systems with Applications* **41(11)**:5252–5266 DOI 10.1016/j.eswa.2014.02.035.
- Canter D, Hammond L, Youngs D, Juszcak P. 2013.** The efficacy of ideographic models for geographical offender profiling. *Journal of Quantitative Criminology* **29(3)**:423–446 DOI 10.1007/s10940-012-9186-6.
- Capozzoli A, Lauro F, Khan I. 2015.** Fault detection analysis using data mining techniques for a cluster of smart office buildings. *Expert Systems with Applications* **42(9)**:4324–4338 DOI 10.1016/j.eswa.2015.01.010.
- Cha SS, Tappert SH, Choi CC. 2010.** A survey of binary similarity and distance measures. *Journal of Systemics, Cybernetics and Informatics* **8(1)**:43–48.

- Chamikara MAP, Galappaththi A, Yapa YPRD, Nawarathna RD, Kodituwakku SR, Gunathilake J, Liyanage LH. 2015.** A crime data analysis framework with geographical information support for intelligence led policing. *PeerJ PrePrints* 3:e1909 DOI 10.7287/peerj.preprints.1529v1.
- Chang YCI, Lin SC. 2004.** Synergy of logistic regression and support vector machine in multiple-class classification. In: *IDEAL 2004. Lecture notes in computer science*, vol. 3177, 132–141.
- Chen H. 1995.** Machine learning for information retrieval: neural. *Journal of the American Society for Information Science* 46(3):194–216 DOI 10.1002/(SICI)1097-4571(199504)46:3<194::AID-ASIA>3.0.CO;2-S.
- Chen H. 2006.** *Intelligence and security informatics for international security*. 1st edition. Vol. 10. Berlin Heidelberg: Springer.
- Chen H, Chung W, Xu JJ, Wang G, Qin Y, Chau M. 2014.** Crime data mining: a general framework and some examples. *Computer* 37(0018-9162):50–56 DOI 10.1109/MC.2004.1297301.
- Chen H, Zeng D, Atabakhsh H, Wyzga W, Schroeder J. 2003.** COPLINK: managing law enforcement data and knowledge. *Communications of the ACM* 46(1):28–34 DOI 10.1145/602421.602441.
- Chikersal P, Poria S, Cambria E. 2015.** SeNTU: sentiment analysis of tweets by combining a rulebased classifier with supervised learning. In: *Proceedings of the international workshop on semanti evaluation (SemEval)*, 647–651.
- Chisum WJ, Turvey B. 2000.** Evidence dynamics: locard’s exchange principle & crime reconstruction. *Journal of Behavioral Profiling* 1(1):1–15.
- Cleary JG, Leonard ET. 1995.** K\*: an instance-based learner using an entropic distance measure. In: *12th international conference on machine learning*, 108–114.
- Douglas JE, Douglas LK. 2006.** Modus operandi and the signature aspects of violent crime. In: *Crime classification manual*. 2nd edition. San Francisco: Jossey-Bass, 19–30.
- Estabrooks A, Jo T, Japkowicz N. 2004.** A multiple resampling method for learning from imbalanced data sets. *Computational Intelligence* 20(1):18–36 DOI 10.1111/j.0824-7935.2004.t01-1-00228.x.
- Fisher D. 1987.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Soybean+%28Small%29>.
- Friedman J, Hastie T, Tibshirani R. 2000.** Additive logistic regression: a statistical view of boosting. *Annals of Statistics* 28(2):337–407.
- Godjevac J. 1997.** *Neuro-fuzzy controllers: design and application*. 1st edition. Lausanne: PPUR presses polytechniques et universitaires romandes.
- Hall M, Frank E, Holmes G, Pfahringer B, Reutemann P, Witten IH. 2009.** The WEKA data mining software. *ACM SIGKDD Explorations Newsletter* 11(1):10 DOI 10.1145/1656274.1656278.
- Hanley JA, McNeil BJ. 1982.** The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology* 143(1):29–36 DOI 10.1148/radiology.143.1.7063747.



- Hempstalk K, Frank E. 2008.** Discriminating against new classes: one-class versus multi-class classification. In: *AI 2008: advances in artificial intelligence: 21st Australasian joint conference on artificial intelligence*. Auckland, 325–336.
- Holdaway S. 1993.** *Issues in sociology: crime and deviance*. Spiral-bound, New edition. Cheltenham: Nelson Thornes Ltd.
- Howell DC.** *Fundamental statistics for the behavioral sciences focuses*. 8th edition. Belmont: Wadsworth, Cengage Learning.
- Hume T. 1994.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Balance+Scale>.
- Iltter N, Guvenir HA. 1998.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Dermatology>.
- Jayawardena AW, Fernando DAK, Zhou MC. 1997.** Comparison of multilayer perceptron and radial basis function networks as tools for flood forecasting. *IAHS Publications-Series of Proceedings and Reports-International Association of Hydrological Sciences* **239**:173–182.
- Julien B. 1990.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Lenses>.
- Keerthi SS, Shevade SK, Bhattacharyya C, Murthy KRK. 2001.** Improvements to Platt's SMO algorithm for SVM classifier design. *Neural Computation* **13**(3):637–649 DOI 10.1162/089976601300014493.
- King RD, Sutton GM. 2013.** High times for hate crimes: explaining the temporal clustering of hate-motivated offending. *Criminology* **51**(4):871–894 DOI 10.1111/1745-9125.12022.
- Koperski K, Han J. 1995.** Discovery of spatial association rules in geographic information databases. In: *Proceeding of the 4th international symposium on spatial databases*, 47–67.
- Kurgan LA, Cios KJ. 2001.** UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/SPECT+Heart>.
- Landwehr N, Hall M, Frank E. 2005.** Logistic model trees. *Machine Learning* **95**(1–2):161–205 DOI 10.1007/s10994-005-0466-3.
- Leclerc B, Proulx J, Beauregard E. 2009.** Examining the modus operandi of sexual offenders against children and its practical implications. *Aggression and Violent Behavior* **14**(1):5–12 DOI 10.1016/j.avb.2008.08.001.
- Lin S, Brown DE. 2006.** An outlier-based data association method for linking criminal incidents. *Decision Support Systems* **41**(3):604–615 DOI 10.1016/j.dss.2004.06.005.
- Mamdani EH, Assilina S. 1975.** An experiment in linguistic synthesis with a fuzzy logic controller. *International Journal of Man-Machine Studies* **7**(1):1–13 DOI 10.1016/S0020-7373(75)80002-2.
- MathWorks. 1994–2015a.** MathWorks. Available at <https://in.mathworks.com/>.
- MathWorks. 1994–2015b.** MathWorks documentation. Available at <http://in.mathworks.com/help/matlab/ref/edit.html>.
- MathWorks. 1994–2015c.** MathWorks fuzzy logic toolbox. Available at <http://in.mathworks.com/help/matlab/ref/edit.html>.

- MathWorks Inc.** 1994–2015. MathWorks. Available at <http://in.mathworks.com/help/fuzzy/foundations-of-fuzzy-logic.html>.
- Michalski RS.** 1980. Learning by being told and learning from examples: an experimental comparison of the two methods of knowledge acquisition in the context of developing an expert system for soybean disease diagnosis. *International Journal of Policy Analysis and Information Systems* 4(2):125–161.
- Mokros A, Alison LJ.** 2002. Is offender profiling possible? Testing the predicted homology of crime scene actions and background characteristics in a sample of rapists. *Legal and Criminological Psychology* 7(1):25–43 DOI 10.1348/135532502168360.
- Oatley G, Ewwart B.** 2011. Data mining and crime analysis. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery* 1(2):147–153 DOI 10.1002/widm.6.
- Palmiotto MJ.** 1988. Crime pattern analysis: an investigative tool. *Critical Issues in Criminal Investigation* 2:59–69.
- Paternoster R, Bachman R.** 2001. *Explaining criminals and crime: essays in contemporary criminological theory*. Los Angeles: Roxbury Publishing Company.
- Pazzani M.** 1991a. UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/Balloons>.
- Pazzani M.** 1991b. The influence of prior knowledge on concept acquisition: experimental and computational results. *Journal of Experimental Psychology: Learning, Memory and Cognition* 17(3):416–432.
- Platt JC.** 1999. Fast training of support vector machines using sequential minimal optimization. In: *Advances in kernel methods*. Cambridge: MIT Press, 185–208.
- Quinlan JR.** 1993. C4.5 programs for machine learning. *Machine Learning* 16(3):235–240.
- Refaeilzadeh P, Tang L, Liu H.** 2009. Cross-validation. In: *Encyclopedia of database systems*. 1st edition. New York: Springer US, 532–538.
- Rish I.** 2001. *An empirical study of the naive Bayes classifier*. Vol. 3. New York, 41–46.
- Sumner M, Frank E, Hall M.** 2005. Speeding up logisti model tree induction. In: *9th European conference on principles and practice of knowledge discovery in databases*, 675–683.
- Sun CT.** 1994. Rule-base structure identification in an adaptive-network-based fuzzy inference system. *IEEE Transactions on Fuzzy Systems* 2(1):64–73 DOI 10.1109/91.273127.
- Takagi T, Sugeno M.** 1985. fuzzy identification of systems and its applications to modeling and control. *IEEE Transactions on Systems, Man and Cybernetics* 1:116–132.
- The ‘Lectric Law Library.** The ‘Lectric Law Library. Available at <http://www.lectlaw.com/files/int20.htm>.
- Thrun S.** 1992. UCI machine learning repository. Available at <https://archive.ics.uci.edu/ml/datasets/MONK's+Problems>.
- Thrun SB, Bala J, Bloedorn E, Bratko I, Cestnik B, Cheng J, De Jong K, Dzeroski S, Fahlman SE, Fisher D, Hamann R, Kaufman K, Keller S, Kononenko I, Kreuziger J, Michalski RS, Mitchell T, Pachowicz PmReich Y, Vafaie H, Van de Welde W, Wenzel W, Wnek J, Zhang J.** 1991. The MONK's problems: performance comparison of different learning algorithms. Technical Report CS-CMU-91-197. Carnegie Mellon University, Pittsburgh.

- Witten IH, Frank E, Trigg L, Hall M, Holmes G, Cunningham SJ. 1999.** Weka: practical machine learning tools and techniques with Java implementations. Available at <http://www.cs.waikato.ac.nz/~ml/publications/1999/99IHW-EF-LT-MH-GH-SJC-Tools-Java.pdf> (accessed 4 June 2016).
- Wolpert DH. 1996.** The lack of a priori distinctions between learning algorithms. *Neural Computation* **8**:1341–1390 DOI [10.1162/neco.1996.8.7.1341](https://doi.org/10.1162/neco.1996.8.7.1341).
- Wu D, Mendel JM. 2007.** Aggregation using the linguistic weighted average and interval type-2 fuzzy sets. *IEEE Transactions on Fuzzy Systems* **15**(6):1145–1161 DOI [10.1109/TFUZZ.2007.896325](https://doi.org/10.1109/TFUZZ.2007.896325).
- Yi X, Rao FY, Bertino E, Bouguettaya A. 2015.** Privacy preserving association rule mining in cloud computing. In: *Proceedings of the 10th ACM symposium on information, computer and communications security*. New York: ACM, 439–450.
- Zadeh LA. 1997.** Toward a theory of fuzzy information granulation and its centrality in human reasoning and fuzzy logic. *Fuzzy Sets and Systems* **90**:117–117.