

CS1800 Situation Analysis

Attached is the **AI Ethics** Situation Analysis for CS1800.

Some things to remember:

- **Think multidimensionally:** cyber policy impacts many different issues, and it is important to identify potential risks and opportunities in your analysis. Consider the strengths and weaknesses of several possible responses and select the optimal response.
- **Engage the scenario:** Assume that the situation which we have provided you is plausible. At the same time, think critically about the information that you have been provided and its origins.
- **Consider interests:** Organizations have a broad and diverse set of interests. How might your decision impact other interests which your organization would like to secure? If you choose one course of action, would a different office at your organization reject your approach? Be sure to be able to justify your response as strongly as possible.
- **Think holistically:** It is important to also consider not just your interests, but *all* parties' interests, including states and non-state actors.

What you are responsible for:

1. A brief, informal oral presentation by your group in section. There are no requirements for using visual materials, nor are there requirements for how many members of your group must speak. Expect to outline your course of action for roughly five, but no more than ten, minutes.
2. Engaging in a brief Q&A from your classmates and your TAs afterwards. We will be assessing the degree to which you have prepared justifications for your course of action.
3. A brief, informal one-page summary of your proposed plan of action to be emailed to your TAs at 11:59 PM the day before your section. Formatting is not especially important – we just want a record of your approach for evaluation purposes. Bullet points are acceptable in this assignment. No bibliography is required.
4. Filling out a peer evaluation after your presentation, during which you will have the opportunity to inform the HTAs about whether all members of your group contributed fairly to your group assignment.

Attached files:

1. A letter requesting your advice on a new tool developed by a Twitter data scientist.
2. A previous (external) attempt at understanding radicalization on social media.
3. An article by the Israeli government describing efforts to profile “lone wolf” extremists.

Disclaimer: While some situation analyses for CS1800 are entirely fictional, others are based on real events. You should think about the “date” of your scenario and discount any events that have occurred in the real world after that date.

You are members of the policy team at Twitter, and have been asked to determine whether and how a new AI tool ought to be used.

April 20, 2020

Dear Colleagues,

We are extremely proud of Twitter’s capacity to expose our users to a wide range of original content. The dark side of this is that our platform may have contributed to the radicalization of individuals who went on to commit violent acts. A study conducted in Israel found that 95% of terrorists used social media before their attacks, and 60% posted information that could have been directly used to predict their attack. Crucially, these “lone wolf” cases are the ones most feared by security services, as they are unable to use traditional methods of identifying extremist networks to predict and thwart attacks. It has typically been difficult to predict lone wolf radicalization: there are a tremendous number of non-radicalized people viewing content that, while extreme, nonetheless meets our terms of service.

A data scientist in our team has proposed a novel solution to this problem of identifying radicalization. This data scientist has developed an artificial neural network (ANN), incorporating Google’s Transformer architecture, to predict which users have been radicalized on our platform and which have simply browsed extreme content. We are then able to use this information to help security services prevent attacks. Previous external efforts have attempted this, but this employee’s method has yielded exceptionally promising results due to the use of new techniques and access to granular non-public information only available internally. To test the network, we provided it with a dataset of 10,000 Twitter users who have had some exposure to an extensive list of far-right and Islamic extremist content, 344 of whom are known to be radicalized. The network correctly predicted the radicalized status of 291 individuals, but incorrectly identified 88 non-radicalized users as radicalized. This is summarized in the table below.

	Predicted as radicalized	Predicted as non-radicalized
Radicalized	291	53
Non-radicalized	88	9,568

While this is the best that has been achieved thus far, it still results in 53 radicalized users not being passed on to the security services, and 88 non-radicalized users unjustly receiving additional scrutiny from their government. Furthermore, employees familiar with the project have voiced privacy

concerns. However, senior officials within the Federal Bureau of Investigation (FBI) in the U.S., the Security Service (MI5) in the U.K., and the Federal Office for the Protection of the Constitution (BfV) in Germany have expressed keen interest in the information resulting from this project, claiming it would be a valuable source of intelligence.

We seek your recommendation on whether to use this network, and if so, how we ought to use it. The data scientist has advised us that while it is possible to decrease the false positives and increase the false negatives or vice versa, it is not possible to decrease both types of errors. Three possible uses have been suggested already: passing positive results to the relevant security services, letting security services query us for our confidence of suspected cases being radicalized, and tweaking our algorithms to make extreme content difficult to find for suspected radicalized users. We are of course open to additional suggestions of how, if at all, to use this tool.

Received March 31, 2017, accepted April 16, 2017, date of publication May 29, 2017, date of current version July 3, 2017.

Digital Object Identifier 10.1109/ACCESS.2017.2706018

Measuring the Radicalisation Risk in Social Networks

RAÚL LARA-CABRERA¹, ANTONIO GONZÁLEZ PARDO¹, KARIM BENOURET²,
NOURA FACI², DJAMAL BENSLIMANE², AND DAVID CAMACHO¹

¹Universidad Autónoma de Madrid, 28049 Madrid, Spain

²Université Claude Bernard Lyon 1, 69622 Villeurbanne, France

Corresponding author: Raúl Lara-Cabrera (raul.lara@uam.es)

This work was supported in part by EphemCH, Spanish Ministry of Economy and Competitiveness, under the European Regional Development Fund FEDER, under Grant TIN2014-56494-C4-4-P and in part by the Justice Programme of the European Union (2014-2020) 723180, RiskTrack, under Grant JUST-2015-JCOO-AG and Grant JUST-2015-JCOO-AG-1.

ABSTRACT Social networks (SNs) have become a powerful tool for the jihadism as they serve as recruitment assets, live forums, psychological warfare, as well as sharing platforms. SNs enable vulnerable individuals to reach radicalized people, hence triggering their own radicalization process. There are many vulnerability factors linked to socio-economic and demographic conditions that make jihadist militants suitable targets for their radicalization. We focus on these vulnerability factors, studying, understanding, and identifying them on the Internet. Here, we present a set of radicalization indicators and a model to assess them using a data set of tweets published by several Islamic State of Iraq and Sham sympathizers. Results show that there is a strong correlation between the values assigned by the model to the indicators.

INDEX TERMS Terrorism, radicalisation, indicator, metric, risk factor, social network analysis.

I. INTRODUCTION

When the 11-S terrorist attack took place in the USA in 2001, the West entered in a new era of continuous danger for its lifestyle as well as its population. This event became a turning point as it started a series of attacks perpetrated by extremists groups on behalf the Islamic State [1].

One of the top priorities of the European Union is to protect the fundamental rights of its citizens as well as to guarantee their safety by fighting all kinds of terrorism. The European Council sets a strategy for Counter-terrorism in 2005 based on four premises: prevention, protection, pursuit and responding [2]. The European Council reworked this strategy in 2014, producing measures and guidelines for the European member states [3].

Although the new jihadist terrorism shares features with other kinds of terrorism it has a distinctive nature that is who radicalise its militants. It is necessary to research the unique traits of jihadist terrorism to recognize the aforementioned peculiarity as well as the different phases a person goes through in order to become radicalised. This understanding may provide relevant information to detect and inhibit radicalisation. The continuous innovation in how the terrorists commit their attacks, the magnitude of the violence perpetrated and the psychological consequences for Western citizens makes the counter-terrorism measures as well as the

prevention of radicalisation critical issues for governments and counter-terrorism institutions [1].

Regarding jihadist radicalisation, there are many vulnerability factors that make their militants suitable targets. These factors are linked to socio-economic and demographic conditions, as stated by the United Nations Office of Drugs and Crime [4]. However, this explanation is very naive (see [5], [6]). Note that there are around 1300 millions Muslims practitioners all over the world who suffer from the same social, economic and political problems as jihadists and, surprisingly, the latter are not as large as it should be expected. Moreover, just a small percentage of these Muslim practitioners are in unison with this fanatic point of view.

In addition to the socio-demographic factors, radicalisation is triggered by feelings, basic needs, emotions as well as by personal life situations and experiences. People usually start their radicalisation by reaching out radical individuals or groups and digging into extremist ideas when they are seeking to fulfil the aforementioned needs. These groups provide social recognition and feel of belonging, which in turn promote the ingress to extremist networks and active membership as well. Speaking about feelings and emotion, guilt, indignation, anger, humiliation, frustration and hatred are the most related ones to jihadist radicalisation [7], [8].

We will focus on these radicalisation factors, studying them on the Internet. Radicalised individuals publish a lot of information in public social media without any security measure, even though there are encryption tools and anonymity software [4] that can be used to "hide" the content of the information. This means that every user of the Social Media can read the majority of the information published by the radicalised individuals. Hence, Social Media is a perfect data source for tracing radicalisation factors as we can access to this plain information.

Social networks such as Tumblr, Instagram, Twitter, Facebook and Youtube have become a powerful propaganda tool for the jihadist cause as they serve as recruitment assets, live forums, psychological warfare and sharing platforms. Many youths have begun to use social networks as a new battleground for the Jihad [9], following the message "any Muslim who tries the Jihad against the enemy by electronic means is considered one way or another a Mijahid" that was published on the Al-Fida and Shumukh al-Islam forum. Moreover, terrorists use Internet to disseminate their propaganda, which is supplemented with justifications, explanations, instructions, slideshows, images and videos, just to cite a few [4]. In addition to social networks, jihadists use other Internet services such as blogs, web pages, forums, emails and peer to peer messaging applications [10]. Therefore, Social Media enable vulnerable individuals to reach radicalised people and even people with the same inquietudes, who may be able to encourage each other's radical ideas, supporting the radicalisation process. Furthermore, these networks promote international communications that bring about feeling of being part of a transnational movement [11], [12].

Self-radicalisation without social interaction is improbable, thus supporting the importance of online connections. Even in cases when the individual seemed to be alone in the process of radicalisation, there were strong influences from people that were already radicalised, or even members of terrorist groups [13]. In some cases, however, the communication with them may be a result of chance [14].

This paper presents some results related to an European Project called Risk-Track¹ [15], whose main goal is the development of a tracking tool based on social media for risk assessment on radicalization. This project is focused on the extraction of radicalization factors on social media and the development of a detection tool. This work corresponds to the first step in the development of a risk assessment tool in Social Networks, and its goal is to define (and validate) the different indicators that later can be used to identify those members of a Social Network with high risk of being radicalised.

The main contributions of this paper are the following: it proposes **five indicators** and their corresponding metrics that can be used to measure the online radicalisation assessment of a given individual using public data from his social networks,

later an experimental evaluation of these indicators, using a public dataset of tweets from several Twitter accounts of pro-ISIS are carried out. Finally, a detailed analysis of the relationships between these metrics are discussed.

The rest of this paper is structured as follows: Section II contains an overview of Social Network Analysis; Section III provides a description of the different radicalisation factors taken into account in this work; A complete description of the proposed model can be found in Section IV and the description of the dataset used is provided in V; Finally, the experimental results obtained in this work are analysed in Section VI and the concluding remarks are explained in Section VII.

II. BACKGROUNDS ON SOCIAL NETWORK ANALYSIS

Social Network Analysis (SNA) is the process of extracting knowledge from Social Network Data. The general process of any SNA platform is composed of the following stages of a pipeline:

- 1) Data Extraction from the Social Networks (SNS).
- 2) Data Preprocessing.
- 3) Data Representation.
- 4) Execution of one, or several, specific algorithm.
- 5) Results analysis.

Any SNA process starts with the extraction of Social Media Data from the specific SN. This first step is required in order to validate the SNA algorithm designed. This process can be skipped if researchers acquire these data from previous own data, or from an existing dataset available on the Internet.

Once the data have been gathered, it is needed to perform some preprocessing. This process is required because typically, these data is not ready for being analysed by the SNA algorithm. For this reason, data need to be processed to ensure that the data can be used by the SNA algorithm. Some of the most typical preprocessing task are the following:

- **Aggregation** is performed when two or more features are combined into a new one.
- **Discretization** is used when a feature with continuous values needs to be represented with discrete values.
- **Normalization** is required when the values of a specific feature needs to be fixed between two upper and lower bounds.
- **Feature Selection** represents the process of subsetting the whole set of features, in such a way only those features relevant to the problem are taken into account.
- **Feature Extraction** consists of transforming current features to generate new ones.
- **Sampling** task consists of extracting a small random subset of instances from the whole data to be processed. The selection process should guarantee that the sample is enough representative of the distribution that governs the data, thereby ensuring that results obtained on the sample are close to ones obtained on the whole dataset.

Once the data has been extracted and preprocessed, researchers have to adapt the data to the best representation for the specific algorithm that will perform the SNA task.

¹<http://risk-track.eu/en/>

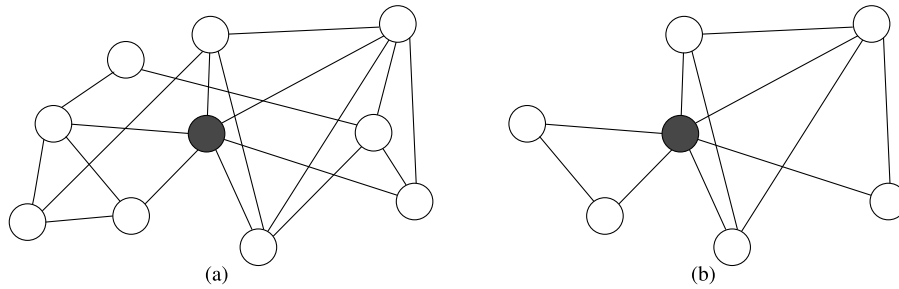


FIGURE 1. a) This figure contains a Social Network composed of 10 users. The different users are represented in the nodes of the graph, whereas the edges represent the different relations between them. b) This figure shows the Ego Network for the coloured node from left SN.

When working with Social Media, the most classical way to represent the data is using a graph. A graph is a structure that can be defined as $G = \{V, E\}$ where V is the set of nodes and E represents the edges of the graph. In the case of social networks, the nodes of the graph represent the different users of the SN, whereas the edges represent the different connections between the users.

Nevertheless, due to the extremely large number of users in Social Networks (SN) [16], researchers try to find reduced representations of the SN to test their algorithms. One of the most extended representations is the well-known Ego Networks [17].

An Ego Network is a social network centered in a specific user called 'Ego'. This network also contains those users connected to this Ego (called 'Alters'), and all the relations between the alters. In any Social Network, there are as many Ego Network as users belong to this SN (i.e. there is one Ego Network per user). This concept is represented in Figure 1, where a small SN composed of 10 users is shown in Figure 1(a), and the corresponding Ego Network for a specific user of this SN can be observed in Figure 1(b).

Ego Networks are used to evaluate the SN and the online social relationships of a specific user. In this sense, Ego Networks are typically used to perform Community Finding tasks in order to find the different communities that compose the contacts of the 'Ego' user [17], [18].

Regarding Community Finding methods, there are three different families of algorithms that can be used. The first family is composed by those algorithms that use the different properties, or topology, of the graph to perform the community finding task. This is the case of the **Clique Percolation Method** (CPM) [19]–[21] that generates the different communities taking into account the connectivity among the different nodes that compose the communities. Other example is **Label Propagation** [22], which uses the topology of the network to propagate several labels that define the different communities.

The second family of algorithms is composed of those algorithms that detect the different communities by performing a hierarchical clustering. In this sense, it is possible to find bottom-up algorithms (i.e. those that start considering that each element of the dataset belongs to a single community,

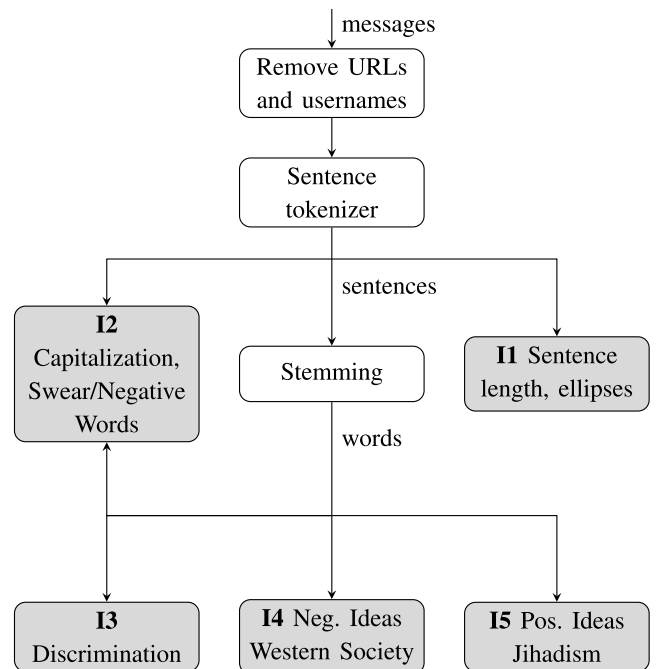


FIGURE 2. Data preprocessing, cleaning and analysis to compute indicators.

and then, iteratively, the different communities are merged according to a specific measure). Examples of bottom-up algorithms are **Clauset-Newman-Moore** [23], [24] (based on the Edge Betweenness algorithm [25], [26]), or the **Walktrap** algorithm [27], [28] that it is based on the random walks that connect the different communities. A different approach within hierarchical clustering are the top-down algorithms. In this case, the algorithms start with only one community that contains all the elements and then, this community is divided taking into account the existing edges or links [29]–[31].

Finally, the third family is composed of probabilistic, or heuristic, approaches, like [32] and [33]. Where the different elements are evaluated by a probabilistic model in order to measure their membership to the different communities.

The proposed work consists of a step backwards to all the aforementioned works. Indeed, our goal is to define the different Online Radicalisation Factors that can be used to

identify those users with high risks of being radicalised. Once these users have been identified, it is possible to build their corresponding Ego Network, and finally analyse the different communities that compose their social network.

III. RADICALISATION FACTORS

As it was mentioned in previous section, the goal of this paper is to define the different indicators that highlight those Social Network users with high risk of being radicalised.

Indicators provide information extracted from the current situation where the individual is involved. As these indicators can launch a process where the selected individual may be investigated, it is extremely important to define correctly these indicators.

In this work, we rely on a list of indicators provided by several expert psychologists on radicalisation that are used to measure the radicalisation level of any individual. Note that there are a huge number of indicators that can be used. For example, Tahir Mahmood has already identified more than 110 indicators extracted from biographies, videos, interviews and information of over 2000 radicalised persons in the world [34].

Nevertheless, as this paper is focused on the usage of Social Media, we will take into account only those indicators that can be measured by the activity of the target users in the Social Networks. We consider **five indicators** (showed in Fig. 2) and group them in two categories, *attitudes* and *beliefs* towards Muslim religion and Western society, and *personality* and *interpersonal relationships*, as grouped by the experts on radicalisation. The former are indicators that are measured by the **content** of the tweets, whereas the later contains those indicators related to the **writing style** specific for each user. We hereafter describe these indicators:

- **Personality related Indicators:**
 - I1** The individual is frustrated.
 - I2** The individual is introverted.
- **Attitudes and beliefs related Indicators:**
 - I3** Perception of discrimination for being Muslim.
 - I4** Expressing negative ideas about Western society.
 - I5** Expressing positive ideas about jihadism.

The first indicator (**I1**) tries to determine whether the user is frustrated. Although this indicator can be determined if the individual is easily irritated, or whether the individual has negative reactions, to measure this indicator in SNS we will take into account some aspects such as the capitalization of the words, or the usage of words with negative content and swearwords.

The goal of the second indicator (**I2**) is to determine whether we are dealing with an introverted user or not. Several studies reveal that introvert users have higher risk of being radicalised. This goal can be computed using the length of the sentences and the usages of ellipses in the tweets.

The third indicator (**I3**) is related to the feeling of being discriminated just because for being Muslim. This perception can be expressed in the content of the tweets, specially if

the individual uses some keywords such as "hate", "sick", "Muslim", etc.

The fourth indicator (**I4**) is their hate to the Western society. This feature can be clearly stated by writing tweets that contain negative ideas about the Western society. Some keywords that can be used are "Western", "hate", "people",

Besides their hate about the Western society, radical people show a deep love for jihadism (**I5**). This feeling can be observed in those tweets or sentences that provide positive ideas about the *mujahideen* (i.e. people engaged in Jihad), or their will to restore the Caliphate. This fifth indicator can be analysed taking into account some keywords like "Islamic State", "Caliphate", or "mujahid".

IV. DESCRIPTION OF THE PROPOSED MODEL

In order to validate the radicalisation indicators presented in Section III, we propose a knowledge extraction model capable of performing a quantitative analysis about written text in social networks. Our proposal uses features from diverse research field such as Natural Language Processing, Data Mining and Statistics. As was shown in Section II, the model consists of various stages through which the texts obtained from Social Networks are processed and analysed.

In the data acquisition stage, the data gathered from the social network (i.e. the posts, tweets and/or status updates) is grouped by the user that posted them. This is a straightforward step that relies solely on the data structure, or storage technology, used when the social network was mined. Next, there is a cleaning step, that is part of the pre-processing stage, to remove the URL's as well as the mention to other users using regular expressions, as this kind of information is out of the scope of the present work and also it could distort the results of the following stages.

During the data representation stage, the messages are divided into sentences because some indicators are based on this unit of language. This task is performed by a sentence tokenizer, that is, a method that divides the text into sentences according to the punctuation marks and new lines it found in the text.

Once the text is divided into sentences it is possible to compute a measurement of the indicator related to the introversion of the individual (**I2**), as this indicator is based on features found at the sentences such as: the use of ellipsis, or the number of words in it (see Section III for a further description). Another feature that might be extracted from sentences is the use of capitalized words associated with frustration (**I2**).

The rest of the indicators (**I3**, **I4**, **I5**) are based on the words used to create the sentence. To measure these indicators it is necessary an additional step to divide each sentence into the different words. As the indicators are computed by counting combined occurrences of some keywords related to the indicator, we decided to broaden the search in two different ways. Firstly, expanding the set of keywords with their synonyms. And secondly, by looking for combined occurrences of the stem of the words. This process is generally known as stem-

ming, and it avoids unwanted situations as not counting as an occurrence the word and its plural form and considering, for instance, ‘fished’ and ‘fish’ as distinct terms. In other words, if two terms share their stem they are considered the same.

Once data is processed and transformed, the indicators can be computed according to the following methodology:

- I1:** Frustration relates to the use of swear words as well as sentences with a negative connotation. To compute these metrics (note that an indicator may have several associated metrics), it is necessary to count the frequencies of swear and negative words and normalize them in a similar way as the other indicators. Furthermore, there is an additional feature that might relate to frustration: capitalization in words, that is, words written with all their letters in upper case. To measure this, the model computes the average number of capitalized words per sentence which is also normalized.
- I2:** Introversion is related to the length of the sentences (usually are short sentences) and the use of ellipses. To measure this indicator the model counts the average sentence length (in number of words) for every user as well as the number of ellipses in his/her tweets, searching sequences of at least three points in the tweet. The value is normalized dividing it between the maximum value obtained to get an indicator within the range [0.0, 1.0].
- I3, I4, I5:** These indicators are related to a set of keywords that express the perception of being discriminated for being Muslim, negative ideas about Western Society and positive ideas about jihadism, respectively. In order to give a numerical value, the model counts how many times there are at least two of the keywords in a sentence (see Table 1 to check the keywords). As with **I2**, values are normalized.

TABLE 1. Initial keywords that are then expanded with synonyms from Wordnet and stemmed.

Indicator	Initial keywords
I1. The individual is frustrated	shit, crap, damn, fuck
I2 Use of words with negative content	hate, guilt, shame, terrible, horrible, bad, fault
I3. Perception of discrimination for being Muslim	Muslim, sick, hate, discrimination, people, racism, religion
I4. Expressing negative ideas about Western society	western, hate, suck, people, west, europe, usa, US, bloody, sick, impure, kuffar, kafir
I5. Expressing positive ideas about jihadism	islamic, state, caliphate, rise, mujahideen, mujahid, help, fight, weapon, gun, weapons

V. EXPERIMENTAL DATASET

It is really difficult to find open datasets related to terrorism, homeland security and radicalisation. In this case, we download a dataset from Kaggle² with over 17000 tweets from

²<https://www.kaggle.com/kzaman/how-isis-uses-twitter>

several Twitter accounts of pro-ISIS since the Paris attacks in November 2015. Data were gathered and processed by a digital agency called Fifth Tribe, and it is released under the *CC0: Public Domain License*.

As stated before, the dataset has about 17000 observations and 8 features, namely: name, user name, description of the account, user’s location where one can put on their Twitter profile, number of followers, number of tweets, date and time when the tweet was posted, and the text itself. Regarding accounts, there are 112 unique user names. There is an interesting observation: there are only 78 unique descriptions among these 112 users, which may suggests that some of the accounts belong to the same person. Using *langid*, a pre-trained language identification tool written in Python [35], we found that most of the tweets (14556 out of 17410) were written in English, followed by 742 and 610 tweets written in Arabic and French, respectively.

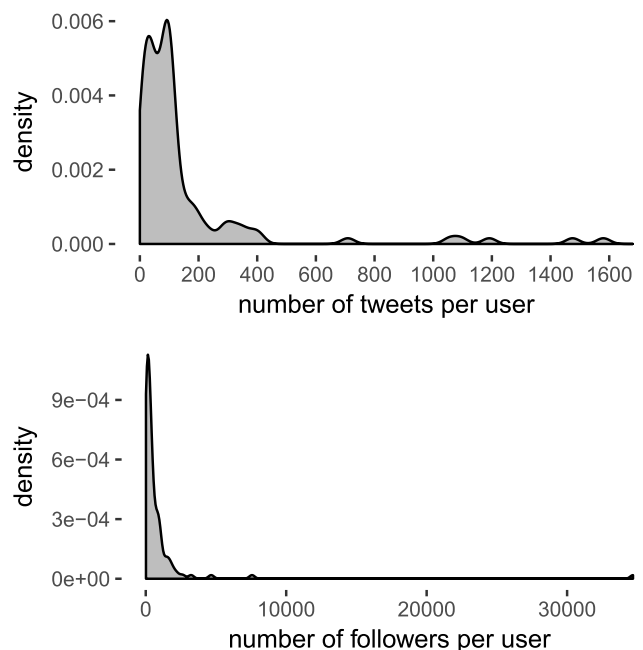


FIGURE 3. Distribution of the number of tweets and followers per user.

As shown in Figure 3, most of the users in the dataset have posted less than 300 tweets. In fact, there are some users with less than 10 published tweets, which is a surprisingly small quantity of information for a user to be tagged as pro-ISIS as the description of the dataset stated. On the other hand, there are some accounts with a high number of tweets that defines themselves as ‘war reporters’. The same thing happens with the number of followers that each user has, whose distribution is heavily left-skewed, with most users having less than 1000 followers and a unique user (@RamiALolah) with over 30000.

According to the timestamps of the tweets in the dataset, they were published between 6th January, 2015 and 13 May, 2016, that is roughly a half and a year. Although the day of the month when the tweets were published are rather chaotic and

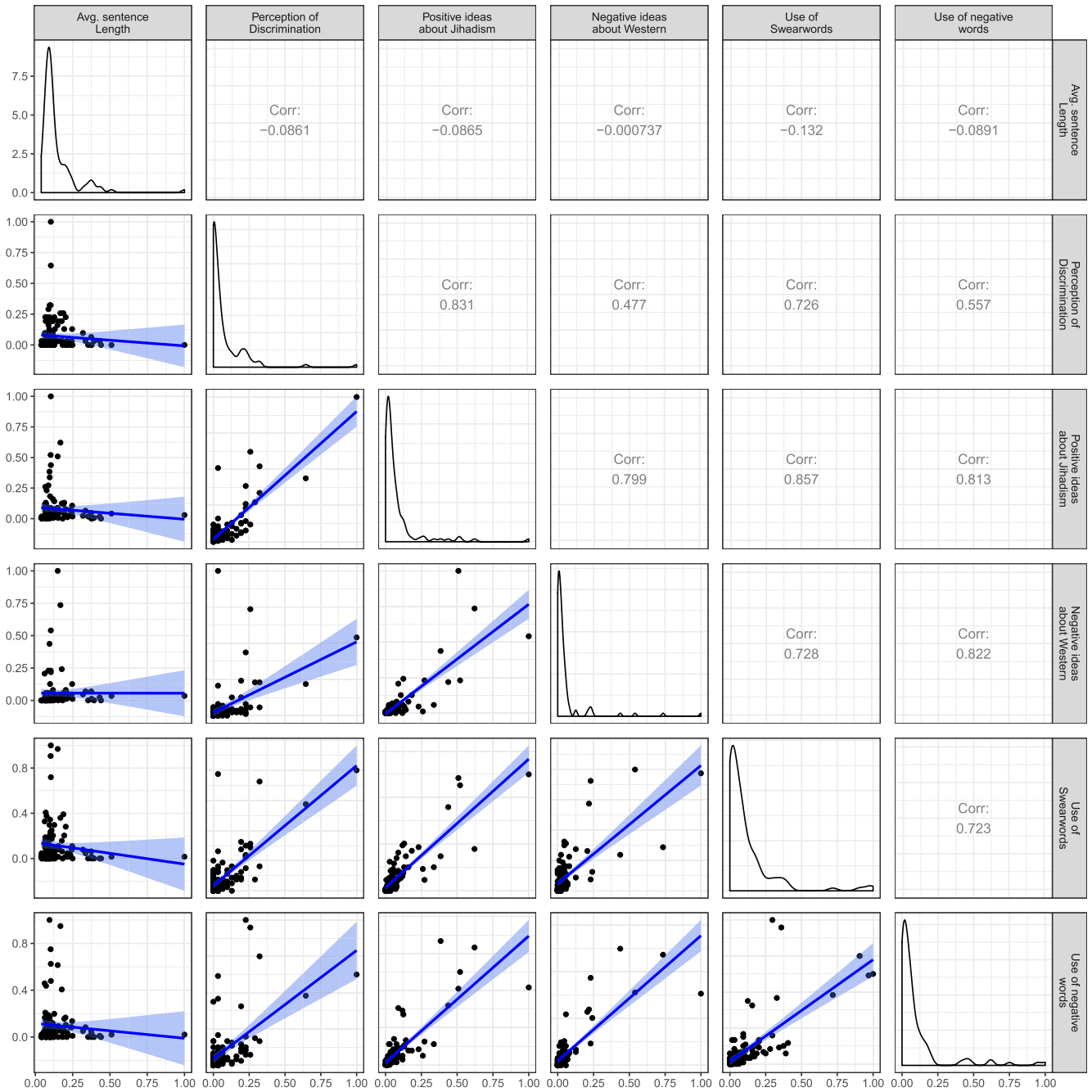


FIGURE 7. Generalized pairs plot of the indicators. The diagonal represents the density plot of the variables. The upper diagonal shows both, the correlation and the Pearson’s correlation coefficient, for each pair of indicators, whereas the lower diagonal represents the scatterplot, and a linear regression model, for those pairs as well (shaded area displays the confidence interval around the model).

words ($\rho = 0.857$, $p\text{-value} < 2.2e-16$) as well as the use of negative words ($\rho = 0.813$, $p\text{-value} < 2.2e-16$). Correspondingly, the linear models computed for the average sentence length against the rest of the metrics do not adjust well to the observations. This situation may happen due to the limit of characters that the social network Twitter sets to the length of the messages, which are limited to 140 characters without counting URL’s. With this limitation in mind it is complicated to extract behaviour information from the average number of words in a sentence, they must be short in order

to express enough information without surpassing the limit of characters. The metric should be useful in other context, though, when there is no limit to the number of characters of a message (Facebook, Tumblr, ...).

VII. CONCLUDING REMARKS

Nowadays, Social Networks (SN) such as Tumblr, Instagram, Twitter, Facebook and Youtube have become a powerful propaganda tool for the jihadist cause as they serve as recruitment assets, live forums, psychological warfare

and sharing platforms. Many youths have begun to use social networks as a new battleground for the Jihad, following the message “any Muslim who tries the Jihad against the enemy by electronic means is considered one way or another a Mujahid” that was published on the Al-Fida and Shumukh al-Islam forum. Moreover, terrorists use Internet to disseminate their propaganda, which is supplemented with justifications, explanations, instructions, slideshows, images and videos, just to name a few. In addition to social networks, jihadists use other Internet services such as blogs, web pages, forums, emails and peer to peer messaging applications. Therefore, Social Media let vulnerable individuals to reach radicalised people and even people with the same inquietudes, who may be able to encourage each other’s radical ideas, supporting the radicalisation process.

This paper defines different indicators that can be used to measure the online radicalisation assessment of a given individual. In this sense, 5 different indicators have been defined related to the attitudes and beliefs towards Muslim religion and Western society and personality and interpersonal relationships. It is important to take into account that this list of indicators can be extended with any other indicator that can be measured using the information extracted from the social networks.

For the experimental phase, we have measured these indicators using a dataset found in Kaggle³ with over 17000 tweets from several Twitter accounts of pro-ISIS since the Paris attacks in November 2015. Data were gathered and processed by a digital agency called Fifth Tribe, and it is released under the *CC0: Public Domain License*.

The first analysis of the indicators reveals that values are distributed on the lower zone of the range of possible values [0.0, 0.1], with some outliers in the mid-range as well as high values in the case of Swearing and Negative Words. This highlights the fact that there is a high number of users with similar values of the indicators and a few users whose indicators are far away from the former.

Also, we have analysed the pairwise correlation between the different indicators. In general, there are strong correlations between the majority of indicators defined in this work: expressing positive ideas about Jihadism, the perception of discrimination, the use of swear words, and the use of negative words. Nevertheless, we have observed a lack of correlation between the average sentence length and the rest of indicators. This situation may happen due to the limit of characters that the social network Twitter sets to the length of the messages, which are limited to 140 characters without counting URL’s. The conclusion that can be drawn from the analysis of these correlations is that it makes sense to measure the indicators with the selected metrics as they share a similar behaviour so people in risk of radicalisation should score high on these metrics; and also, that the dataset used is coherent.

This paper supposes an important step in the fight against radicalisation because it defines several indicators that can

be used to measure the risk of radicalisation of a given individual. Once this individual has been identified, different actions can be performed in order to avoid this individual become a Jihadist. From the computational science point of view, one of this actions could be the analysis of his/her online friendships to identify those communities that influence the user to become a *mujahideen*.

ACKNOWLEDGEMENTS

The contents of this publication are the sole responsibility of their authors and can in no way be taken to reflect the views of the European Commission.

REFERENCES

- [1] D. Garriga Guitart, *Yihad: ¿qué es?*. Barcelona, Spain: Comanegra, 2015.
- [2] Council of the European Union, *The European Union Strategy for Combating Radicalisation and Recruitment to Terrorism*. Brussels, Belgium: European Parliament, 2005.
- [3] European Commission. (2016). *Radicalisation*. [Online]. Available: http://ec.europa.eu/dgs/home-affairs/what-we-do/policies/crisis-and-terrorism/radicalisation/index_en.htm
- [4] United Nations Office On Drugs and Crime, “The use of the Internet for terrorist purposes,” United Nations Office Drugs Crime, Vienna, Austria, Tech. Rep., 2012.
- [5] L. De la Corte Ibáñez and J. Jordán, *La yihad terrorista*. Madrid, Spain: Síntesis, 2007.
- [6] L. De la Corte Ibáñez, *La Lógica del Terrorismo*. Madrid, Spain: Alianza Editorial, 2014.
- [7] S. Atran, *Talking to the Enemy: Religion, Brotherhood, and the (Un)Making of Terrorists*. New York, NY, USA: Harper Collins, 2010.
- [8] A. Speckhard, *Talking to Terrorists: Understanding the Psycho-Social Motivations of Militant Jihadi Terrorists, Mass Hostage Takers, Suicide Bombers & Mart Paperback*. McLean, VA, USA: Advances Press, 2012.
- [9] E. L. Illaro, “Terrorismo islamista en las redes—La yihad electrónica,” Inst. Español Estudios Estratégicos, Madrid, Spain, Tech. Rep. 100/2015, 2015.
- [10] F. J. Cilluffo, S. L. Cardash, and A. J. Whitehead, “Radicalization: Behind bars and beyond borders,” *Brown J. World Affairs*, vol. 13, no. 2, pp. 113–122, 2007.
- [11] S. Ulph, “A guide to jihad on the Web,” *Terrorism Focus*, vol. 2, no. 7, 2005.
- [12] M. R. Torres, *El ECO Del Terror: Ideología y Propaganda en el Terrorismo Yihadista*, Madrid, Spain: Plaza, 2009.
- [13] M. T. Buezo, “El lobo solitario como elemento emergente y evolución táctica del terrorismo yihadista,” *Inteligencia y Seguridad: Revista de Análisis y Prospectiva*, vol. 2013, no. 14, pp. 117–150, 2013.
- [14] M. Sageman, *Understanding Terror Networks*. Philadelphia, PA, USA: Univ. Pennsylvania Press, 2004.
- [15] C. Camacho, A. Gonzalez-Pardo, A. Ortigosa, I. Gilperez-Lopez, and C. Urruela, “Risktrack: A new approach for risk assessment on radicalisation based on social media data,” in *Proc. Workshop Affect. Comput. Context Awareness Ambient Intell. (AfCAI)*, vol. 1794. 2016, pp. 1–10.
- [16] K. Musiał and P. Kazienko, “Social networks on the Internet,” *World Wide Web*, vol. 16, no. 1, pp. 31–72, 2012.
- [17] A. Gonzalez-Pardo, J. J. Jung, and D. Camacho, “ACO-based clustering for ego network analysis,” *Future Generat. Comput. Syst.*, vol. 66, pp. 160–170, Jan. 2017.
- [18] J. Xie, S. Kelley, and B. K. Szymanski, “Overlapping community detection in networks: The state-of-the-art and comparative study,” *ACM Comput. Surv.*, vol. 45, no. 4, pp. 43:1–43:35, 2013.
- [19] A. Lancichinetti and S. Fortunato, “Community detection algorithms: A comparative analysis,” *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 80, pp. 1–11, Nov. 2009.
- [20] S. Gregory, “Finding overlapping communities in networks by label propagation,” *New J. Phys.*, vol. 12, no. 10, pp. 2–27, 2010.
- [21] S. Günemann, B. Boden, I. Färber, and T. Seidl, “Efficient mining of combined subspace and subgraph clusters in graphs with feature vectors,” in *Proc. Pacific-Asia Conf. Knowl. Discovery Data Mining*, 2013, pp. 261–275.

³<https://www.kaggle.com/kzaman/how-isis-uses-twitter>

- [22] U. N. Raghavan, R. Albert, and S. Kumara, "Near linear time algorithm to detect community structures in large-scale networks," *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 76, Sep. 2007, Art. no. 036106.
- [23] Y. Ding, "Community detection: Topological vs. topical," *J. Informetrics*, vol. 5, no. 4, pp. 498–514, 2011.
- [24] S. Sobolevsky, R. Campari, A. Belyi, and C. Ratti, "General optimization technique for high-quality community detection in complex networks," *Phys. Rev. E, Stat. Phys. Plasmas Fluids Relat. Interdiscip. Top.*, vol. 90, no. 1, 2014, Art. no. 012811.
- [25] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi, "Defining and identifying communities in networks," *Proc. Nat. Acad. Sci. USA*, vol. 101, no. 9, pp. 2658–2663, 2004.
- [26] S. Moon, J.-G. Lee, and M. Kang, "Scalable community detection from networks by computing edge betweenness on MapReduce," in *Proc. Int. Conf. Big Data Smart Comput. (BIGCOMP)*, 2014, pp. 145–148.
- [27] P. Pons and M. Latapy, "Computing communities in large networks using random walks," *J. Graph Algorithms Appl.*, vol. 10, no. 2, pp. 191–218, 2006.
- [28] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre, "Fast unfolding of communities in large networks," *J. Stat. Mech., Theory Experim.*, vol. 2008, no. 10, p. P10008, 2008.
- [29] R. Balasubramanian and W. W. Cohen, "Block-LDA: Jointly modeling entity-annotated text and entity-entity links," in *Handbook of Mixed Membership Models and Their Applications*. London, U.K.: Chapman & Hall, 2014, pp. 255–273.
- [30] Y.-Y. Ahn, J. P. Bagrow, and S. Lehmann, "Link communities reveal multiscale complexity in networks," *Nature*, vol. 466, no. 7307, pp. 761–764, 2010.
- [31] T. S. Evans and R. Lambiotte, "Line graphs of weighted networks for overlapping communities," *Eur. Phys. J. B*, vol. 77, no. 2, pp. 265–272, 2010.
- [32] J. Yang and J. Leskovec, "Overlapping community detection at scale: A nonnegative matrix factorization approach," in *Proc. 6th ACM Int. Conf. Web Search Data Mining (WSDM)*, 2013, pp. 587–596.
- [33] J. McAuley and J. Leskovec, "Discovering social circles in ego networks," *ACM Trans. Knowl. Discov. Data*, vol. 8, no. 1, Feb. 2014, Art. no. 4.
- [34] T. Mahmood. (Feb. 2016). *Indicators of Radicalisation—Quality Control Initiative*. [Online]. Available: <https://www.linkedin.com/pulse/indicators-radicalisation-quality-control-initiative-tahir-mahmood>
- [35] M. Lui and T. Baldwin, "Cross-domain feature selection for language identification," in *Proc. 5th Int. Joint Conf. Natural Lang. Process.*, 2011, pp. 553–561.
- [36] G. Salton and J. Michael, *Introduction to Modern Information Retrieval*. New York, NY, USA: McGraw-Hill, 1983, pp. 24–51.



computing, and crowdsourcing.

KARIM BENOURET was a Teaching and Research Assistant with Télécom Saint-Étienne, and a Post-Doctoral Researcher with Inria Nancy–Grand Est. He is currently an Associate Professor with Université Claude Bernard Lyon 1. He is also a member of the Lyon Research Center for Images and Intelligent Information Systems associated with the French National Center for Scientific Research. His research interests include preference queries, recommender systems, services



NOURA FACI has been an Associate Professor with University Lyon 1, France, since 2008. She has authored several papers in high-quality journals and conferences and actively contributes to the IEEE TSC, IEEE IC, and *Computer* journal's review process, as well. Her current research interests are service computing, social computing, and business process management.



tation systems, Web services, ontologies, and databases.

DJAMAL BENSLIMANE is currently a Full Professor of computer sciences with Lyon 1 University and a member of the Service Oriented Computing Research Team, Lyon Research Center for Images and Information Systems, Lyon, France. He has authored papers in well-known journals including IEEE TKDE, ACM TOIT, SIGMOD Record, IEEE TSC, IEEE INTERNET COMPUTING, WWW Journal, and DPDB Journal. His research interests include distributed informa-



RAÚL LARA-CABRERA received the M.Sc. and Ph.D. degrees in computer science from the University of Málaga, Spain, in 2013 and 2015, respectively. He is currently a Research Fellow with the Department of Ingeniería Informática, Universidad Autónoma de Madrid, Spain. His main research areas involve computational intelligence, video games, and complex systems.



PSO, and SWARM intelligence), multi-agent systems, and machine learning techniques. The application domains for his research are constraint satisfaction problems, complex graph-based problems, optimization problems, and video games.

ANTONIO GONZÁLEZ PARDO received the B.Sc. degree in computer science from the Universidad Carlos III de Madrid in 2009, the M.Sc. degree in computer science from the Universidad Autónoma de Madrid in 2011, and the Ph.D. degree in computer science from the Universidad Autónoma de Madrid in 2014. He is currently a Lecturer with the Universidad Autónoma de Madrid. His main research interests are related to computational intelligence (genetic algorithms,



data mining (clustering), evolutionary computation (GA & GP), multi-agent systems and swarm intelligence (Ant colonies), automated planning and machine learning, or video games among others.

DAVID CAMACHO received the B.S. in physics from the Universidad Complutense de Madrid in 1994 and the Ph.D. degree in computer science from the Universidad Carlos III de Madrid in 2001. He is currently an Associate Professor with the Computer Science Department, Universidad Autónoma de Madrid, Spain, and the Head of the Applied Intelligence and Data Analysis Group. He has authored over 200 journals, books, and conference papers. His research interests include



Ministry of Public Security

News

Study: Terrorists post info on social media before attacking

Subject: [International Homeland Security Forum](#) [Crime and society](#)

Publish Date: 12.06.2018

The study on 'lone wolf' terrorists will be presented at the International Homeland Security Forum this week



Findings from a new study being presented at the [International Homeland Security Forum](#), led by the Minister of Public Security Gilad Erdan and the U.S. Secretary of Homeland Security Kirstjen Nielsen, shed light on the “lone wolf” phenomenon and the psychological and sociological profile of attackers in the recent wave of terrorism over the past year and a half.

The study, conducted by Professors Ariel Merari and Boaz Ganor of the International Institute for Counter-Terrorism (ICT) at IDC Herzliya, in partnership with the Israel Ministry of Public Security, focused on independent terrorists (lone wolves) in the surge of terrorism that Israel experienced from October 2015 to December 2017. These are terrorists that acted alone or with accomplices, but with no operational support from a terrorist organization.

The study is based on a database of 700 attackers who took part in 560 attacks.

This is a first of its kind study that includes interviews with 45 lone wolf terrorists in prison, utilizing a unique method to create profiles of the terrorists and their motivations for the attacks, with the hope of formulating techniques to prevent attacks.

One of the alarming figures found was that at least 95 attackers made use of social media prior to carrying out attacks, with 60 of them even publishing posts indicating their intent to carry out attacks.

The study also examined the “success rate” of the attacks (defined as injuring at least one victim) and found that during the first stage of the surge (October 2015-March 2016), over 50% of attacks were successful, while during the second stage (April 2016-December 2017), only 26% of attacks were successful.

The study looked at the various motives for committing terrorist attacks, with in-depth interviews conducted on 45 terrorists, and found that a combination of motives and factors were behind the attacks, including psychological factors, ideological motives, personal factors and trigger events (copycat attacks, geopolitical events and traumatic events) – all augmented by incitement. Ideological motivation (nationalistic and religious) was found to have had an effect on 60% of the terrorists in the sample group – 28% of the men and 11% of the women.

It is important to note that two-thirds of the attackers in the sample group suffered from mental disorders, psychosis or suicidal tendencies. A large percentage of the sample group was suicidal, with 54% saying they would have preferred to die in the attack.

The study found that of the 700 lone wolf attackers, 85% were men and 15% were women, with the average age being 22. Among the attackers, 77% were residents of Judea and Samaria, and 17% were residents of East Jerusalem.

Among the sample group interviewed, there was a particularly high tendency of familial problems among the female attackers. Additionally, 15% of the sample group was found to be illiterate males.

More on the subject

[More about the International Homeland Security Forum](#)

This page was last updated on 12.06.2018