

Predicting online extremism, content adopters, and interaction reciprocity

Emilio Ferrara*, Wen-Qiang Wang*, Onur Varol†, Alessandro Flammini† and Aram Galstyan*

*Information Sciences Institute, University of Southern California, Marina del Rey, CA 90292, USA

†School of Informatics and Computing, Indiana University, Bloomington, IN 47401, USA

emiliofe@usc.edu, wenqianw@usc.edu, ovarol@indiana.edu, aflammin@indiana.edu, galstyan@isi.edu

Abstract—We present a machine learning framework that leverages a mixture of metadata, network, and temporal features to detect extremist users, and predict content adopters and interaction reciprocity in social media. We exploit a unique dataset containing millions of tweets generated by more than 25 thousand users who have been manually identified, reported, and suspended by Twitter due to their involvement with extremist campaigns. We also leverage millions of tweets generated by a random sample of 25 thousand regular users who were exposed to, or consumed, extremist content. We carry out three forecasting tasks, (i) to detect extremist users, (ii) to estimate whether regular users will adopt extremist content, and finally (iii) to predict whether users will reciprocate contacts initiated by extremists. All forecasting tasks are set up in two scenarios: a *post hoc* (time independent) prediction task on aggregated data, and a simulated real-time prediction task. The performance of our framework is extremely promising, yielding in the different forecasting scenarios up to 93% AUC for extremist user detection, up to 80% AUC for content adoption prediction, and finally up to 72% AUC for interaction reciprocity forecasting. We conclude by providing a thorough feature analysis that helps determine which are the emerging signals that provide predictive power in different scenarios.

1. Introduction

Researchers are devoting increasing attention to the issues related to online extremism, terrorist propaganda and radicalization campaigns [1], [2]. Social media play a central role in these endeavors, as increasing evidence from social science research suggests [3], [4]. For example, a widespread consensus on the relationship between social media usage and the rise of extremist groups like the Islamic State of Iraq and al-Sham (viz. ISIS) has emerged among policymakers and security experts [5], [6], [7]. ISIS' success in increasing its roster to thousands of members has been related in part to a savvy use of social media for propaganda and recruitment purposes. One reason is that, until recently, social media platform like Twitter provided a public venue where single individuals, interest groups, or organizations, were given the ability to carry out extremist discussions and terrorist recruitment, without any form of restrictions, and with the possibility of gathering audiences of potentially millions. Only recently, some mechanisms have been put

into place, based on manual reporting, to limit these abusive forms of communications. Based on this evidence, we argue in favor of developing computational tools capable of effectively analyzing massive social data streams, to detect extremist users, to predict who will become involved in interactions with radicalized users, and finally to determine who is likely to consume extremist content. The goal of this article is to address the three questions above by proposing a computational framework for detection and prediction of extremism in social media. We tapped into Twitter to obtain a relevant dataset, leveraged expert crowd-sourcing for annotation purposes, and then designed, trained and tested the performance of our prediction system in static and simulated real-time forecasts, as detailed below.

Contributions of this work

The main contributions of our work can be summarized as:

- We formalize three different forecasting tasks related to online extremism, namely the detection of extremist users, the prediction of adoption of extremist content, and the forecasting of interaction reciprocity between regular users and extremists.
- We propose a machine prediction framework that analyzes social media data and generates features across multiple dimensions, including user metadata, network statistics, and temporal patterns of activity, to perform the three forecasting tasks above.
- We leverage an unprecedented dataset that contains over 3 millions tweets generated by over 25 thousand extremist accounts, who have been manually identified, reported, and suspended by Twitter. We also use around 30 million tweets generated by a random sample of 25 thousand regular users who were exposed to, or consumed, extremist content.
- For each forecasting task, we design two variants: a post-hoc (time independent) prediction task performed on aggregated data, and a simulated real-time forecast where the learning models are trained as if data were available up to a certain point in time, and the system must generate predictions on the future.
- We conclude our analysis by studying the predictive power of the different features employed for prediction, to determine their role in the three forecasts.

2. Data and Preliminary Analysis

In this section we describe our dataset, the curation strategy yielding the annotations, and some preliminary analysis.

2.1. Sample selection and curation

In this work we rely on data and labels constructed by using a procedure of manual curation and expert verification. We retrieved on a public Website a list of over 25 thousands Twitter accounts whose activity was labeled as supportive of the Islamic State by the crowd-sourcing initiative called *Lucky Troll Club*. The goal of this project was to leverage annotators with expertise in Arabic languages to identify ISIS accounts and report them to Twitter. Twitter’s anti-abuse team manually verifies all suspension requests, and grants some based on the active violation of Twitter’s Terms of Service policy against terrorist- or extremist-related activity. Here we focus on the 25,538 accounts that have been all suspended between January and June 2015 by Twitter as a consequence of evidence of activity supporting the Islamic State group. For each account, we also have at our disposal information about the suspension date, and the number of followers of that user as of the suspension date.

2.2. Twitter data collection

The next step of our study consisted in collecting data related to the activity of the 25,538 ISIS supporters on Twitter. To this purpose, we leveraged the Twitter *gardenhose* data source (roughly 10% of the Twitter stream) collected by Indiana University [8]. We decided to collect not only the tweets generated by these accounts prior to their suspension, but also to build a dataset of their targets. In particular, we are concerned with accounts unrelated to ISIS with whom the ISIS supporters tried to establish some forms of interaction. We therefore constructed the following two datasets:

ISIS accounts: this dataset contains 3,395,901 tweets generated in the time interval January-June 2015 by the 25,538 accounts identified by Twitter as supporters of ISIS. This is a significant portion of all the accounts suspended by Twitter in relation to ISIS.¹

Users exposed to ISIS: this dataset contains 29,193,267 tweets generated during January-June 2015 by a set of 25 thousand users randomly sampled among the larger set of users that constitute ISIS accounts’ followers. This set is by choice of equal size to the former one, to avoid introducing class imbalance issues.

For prediction purposes, we will use as positive and negative labels the *ISIS accounts* group and the accounts in the *users exposed to ISIS*, respectively.

1. The Guardian recently reported that between April 2015 and February 2016, Twitter’s anti-terror task force suspended about 125,000 accounts linked to ISIS extremists: <http://www.theguardian.com/technology/2016/feb/05/twitter-deletes-isis-accounts-terrorism-online>

3. Methodology

In this section we discuss the learning models and the features adopted by our framework. The complete prediction pipeline (learning models, cross validation, feature selection, and performance evaluation) is developed using Python and the scikit-learn library [9].

3.1. Learning models

We adopt two off-the-shelf learning models as a proof of concept for the three classification tasks that we will discuss later (see §4): *Logistic Regression* and *Random Forests*.

Logistic Regression: The first implemented algorithm is a simple Logistic Regression (LR) with LASSO regularization. The advantage of this approach is its scalability, which makes it very effective to (possibly real-time) classification tasks on large datasets. The only parameter to tune is the loss function C . We expect that LR will provide the baseline classification and prediction performance.

Random Forests: We also use a state-of-the-art implementation of *Random Forests* (RF) [10]. The vectors fed into the learning models represent each user’s features. *Random Forests* are trained using 100 estimators and adopting the Gini coefficient to measure the quality of splits. Optimal parameters setting is obtained via cross validation (see 3.1.1).

Note that the goal of this work is not to provide new machine learning techniques, but to illustrate that existing methods can provide promising results. We also explored additional learning models (e.g., SVM, Stochastic Gradient Descent, etc.), which provide comparable prediction performance but are less computationally efficient and scalable.

3.1.1. Cross validation. The results of our performance evaluation (see §4) are all obtained via k -fold cross validation. We adopt $k = 5$ folds, and therefore use 80% of data for training, and the remainder 20% for testing purpose, averaging performance scores across the 5 folds. We also use 5-fold cross validation to optimize the parameters of the two learning algorithms (LR and RF), by means of an exhaustive cross-validated grid search on the hyperparameter space.

3.1.2. Evaluation scores. We benchmark the performance of our system by using four standard prediction quality measures, namely Precision, Recall, F1 (harmonic mean of Precision and Recall), and AUC—short for Area Under the Receiver Operating Characteristic (ROC) curve [11].

3.2. Feature engineering and feature selection

We manually crafted a set of 52 features belonging to three different classes: user metadata, timing features, and network statistics, as detailed below.

User metadata and activity features: User metadata have been proved pivotal to model classes of users in social media [12], [13]. We build user-based features leveraging the metadata provided by the Twitter API related to the

TABLE 1. LIST OF 52 FEATURES EXTRACTED BY OUR FRAMEWORK

User metadata & activity	Number of followers
	Number of friends (i.e., followees)
	Number of posted tweets
	Number of favorite tweets
	Ratio of retweets / tweets
	Ratio of mentions / tweets
	Avg number of hashtags
	(avg, var) number of retweets
	Avg. number of mentions
	Avg. number of mentions (excluding retweets)
	Number of URLs in profile description
Timing	(avg, std, min, max, proportion) URLs in tweets
	Length of username
	(avg, var) number of tweets per day
Netw. stats	(avg, std, min, max) interval between two consecutive tweets
	(avg, std, min, max) interval between two consecutive retweets
	(avg, std, min, max) interval between two consecutive mentions
	(avg, std, min, max) distribution of retweeters' number of followers
	(avg, std, min, max) distribution of retweeters' number of friends
	(avg, std, min, max) distribution of mentioners' number of followers
	(avg, std, min, max) distribution of mentioners' number of friends
	(avg, std, min, max) number of retweets of the tweets by others

author of each tweet, as well as the source of each retweet. User features include the number of tweets, followers and friends associated to each users, the frequency of adoption of hashtags, mentions, and URLs, and finally some profile descriptors. In total, 18 user metadata and activity features are computed (cf. Table 1).

Timing features: Important insights may be concealed by the temporal dimension of content production and consumption, as illustrated by recent work [14], [15]. A basic timing feature is the average number of tweets posted per day. Other timing features include statistics (average, standard deviation, minimum, maximum) of the intervals between two consecutive events, e.g., two tweets, retweets, or mentions. Our framework generates 14 timing features (cf. Table 1).

Network statistics: Twitter content spreads from person to person via retweets and mentions. We expect that the emerging network structure carries important information to characterize different types of communication. Prior work shows that using network features significantly helps prediction tasks like social bot detection [8], [13], [15], and campaign detection [16], [17]. Our framework focuses on two types of networks: (i) retweet, and (ii) mention networks. Users are nodes of such networks, and retweets or mentions are directed links between pairs of users. For each user, our framework computes the distribution of followers and friends of all users who retweet and mention that user, and extracts some descriptive statistics (average, standard deviation, minimum, maximum) of these distributions. Our system builds 20 network statistics features (cf. Table 1).

3.2.1. Greedy feature selection. Our framework generates a set F of $|F| = 52$ features. In our type of prediction tasks, some features exhibit more predictive power than others: temporal dependencies introduce strong correlations among some features, thus some possible redundancy. Among the different existing ways to select the most relevant features

for the prediction task at hand [18], in the interest of computational efficiency, we adopted a simple greedy forward feature selection method, as follows: (i) initialize the set of selected features $S = \emptyset$; (ii) for each feature $f \in F - S$, consider the union set $U = S \cup f$; (iii) train the classifier using the features in U ; (iv) test the average performance of the classifier trained on this set; (v) add to S the feature that provides the best performance; (vi) repeat (ii)-(v) as long as a significant performance increase is yield.

4. Experiments

In the following, we formalize three different prediction problems related to online extremism:

Task I (T1): Detection of extremist supporters. The first task that our system will face is a binary classification aimed to detect ISIS accounts and separate them from those of regular users. The problem is to test whether any predictive signal is present in the set of features we designed to characterize social media activity related to extremism, and serves as a yardstick for the next two more complex problems.

Task II (T2): Predicting extremist content adoption. The set of 25 thousand users we randomly sampled among followers of ISIS accounts can be leveraged to perform the prediction of extremist content adoption. We define as a positive instance of adoption in this context when a regular user retweets some content s/he is exposed to that is generated by an ISIS account.

Task III (T3): Predicting interactions with extremists. The third task presents likely the most difficult challenge: predicting whether a regular user will engage into interactions with extremists. A positive instance of interaction is represented by a regular user replying to a contact initiated by an ISIS account.

Static versus real-time predictions. For each of the three prediction tasks described above, we identified two modalities, namely a static (time independent) and a simulated real-time prediction. In the former scenario, a static prediction ignores temporal dependencies in that the system aggregates all data available across the entire time range (January-June 2015), and then performs training and testing using the 5-fold cross validation strategy by randomly splitting datapoints into the 5 folds and averaging the prediction performance across folds. In the latter scenario, a real-time prediction is simulated in which data are processed for training and testing purposes by respecting the timeline of content availability: for example, the system can exploit the first month of available data (January 2015) for training, and then producing predictions for the remainder 5 months (Feb-Jun 2015), for which performance is tested.

The performance of our framework in the three tasks, each with the two prediction modalities, is discussed in the following. The section concludes with the analysis of feature predictive power (see §4.4).

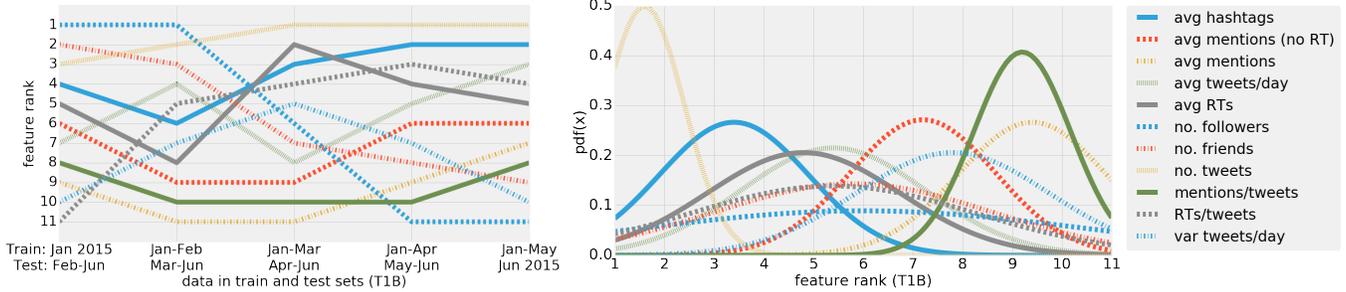


Figure 1. T1B: Feature selection analysis and feature rank distribution (top 11 features)

4.1. T1: Detection of extremist supporters

In the following we discuss the static (T1A) and real-time (T1B) scenarios for the first prediction task, namely detecting extremist accounts on Twitter.

4.1.1. T1A: Time-independent detection. The detection of extremist user accounts is the most natural task to start the performance evaluation of our framework. Our analysis aims at verifying that the 52 features we carefully hand-crafted indeed carry some predictive signal useful to separate extremist users from regular ones. The dataset at hand contains two roughly equal-sized classes (about 25 thousand instances each), where ISIS accounts are labeled as positive instances, and regular users as negative ones. Each instance is characterized by a 52-dimensional vector, and positive and negative examples are fed to the two learning models (LR and RF). The first task, in short T1A, is agnostic of time dependencies: data are aggregated throughout the entire 6 months period (January-June 2015) and training/testing is performed in a traditional 5-fold cross-validated fashion (cf. §3.1.1). Table 2 summarizes the performance of LR and RF according to the four quality measures described above (cf. 3.1.2): Both models perform well, with *Random Forests* achieving an accuracy above 87% as measured by AUC. These results are encouraging and demonstrate that simple off-the-shelf models can yield good performance in T1A.

TABLE 2. EXTREMISTS DETECTION (T1A)

	Precision	Recall	F1	AUC
Logistic Regression	0.778	0.506	0.599	0.756
Random Forests	0.855	0.893	0.874	0.871

4.1.2. T1B: Simulated real-time detection. A more complex variant of this prediction task is by taking into account the temporal dimension. Our prior work has demonstrated that accounting for temporal dependencies is very valuable in social media prediction tasks and significantly improves prediction performance [17]: therefore we expect that the performance of our framework in a simulated real-time prediction task will exceed that of the static scenario.

In this simulated real-time prediction task, T1B, we divide the available data into temporal slices used separately for training and prediction purposes. Table 3 reports five columns, each of which defines a scenario where one or

more months of data are aggregated for training, and the rest is used for prediction and performance evaluation. For example, in the first column, the learning models are trained on data from January 2015, and the prediction are performed and evaluated on future data in the interval February-June 2015.

Random Forests greatly benefits from accounting for temporal dependencies in the data, and the prediction performance as measured by AUC ranges between 83.8% (with just one month of training data) to an excellent 93.2% (with five months of training data). Fig. 1(left) illustrates the ranking of the top 11 features identified by feature selection, as a function of the number of months of data in the training set. Fig. 1(right) displays the distributions of the rankings of each feature across the 5 different temporal slices. For the extremist users detection task, the most predictive features are (1) number of tweets, (2) average number of hashtags, and (3) average number of retweets. One hypothesis is that extremist users are more active than average users, and therefore exhibit distinctive patterns related to volume and frequency of activity.

TABLE 3. REAL-TIME EXTREMISTS DETECTION (T1B)

Training:	Jan	Jan-Feb	Jan-Mar	Jan-Apr	Jan-May
Testing:	Feb-Jun	Mar-Jun	Apr-Jun	May-Jun	Jun
AUC (LR)	0.743	0.753	0.655	0.612	0.602
Precision (LR)	0.476	0.532	0.792	0.816	0.796
Recall (LR)	0.629	0.675	0.377	0.289	0.275
F1 (LR)	0.542	0.595	0.511	0.427	0.409
AUC (RF)	0.838	0.858	0.791	0.942	0.932
Precision (RF)	0.984	0.922	0.868	0.931	0.910
Recall (RF)	0.679	0.733	0.649	0.957	0.959
F1 (RF)	0.804	0.817	0.743	0.944	0.934

4.2. T2: Predicting extremist content adoption

The second task, namely predicting the adoption of extremist content by regular users, is discussed in the static (T2A) and real-time (T2B) scenarios in the following.

4.2.1. T2A: Time-independent prediction. The first instance of T2 is again on the time-aggregated datasets spanning January-June 2015. Predicting content adoption is a known challenging task, and a wealth of literature has explored the factors behind online information contagion [19].

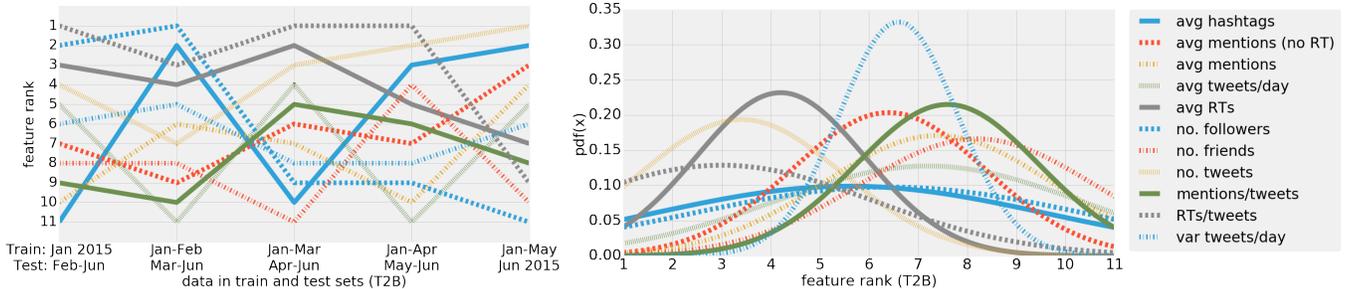


Figure 2. T2B: Feature selection analysis and feature rank distribution (top 11 features)

In this scenario, we aim to predict whether a regular user will retweet a content produced by an ISIS account. Positive instances are represented by users who retweeted at least one ISIS tweet in the aggregated time period, while negative ones are all users exposed to such tweets who did not retweet any of them. Table 4 summarizes the performance of our models: *Random Forests* emerges again as the best performer, although the gap with *Logistic Regression* is narrow, and the latter provides significantly better Recall. Overall, T2A appears clearly more challenging than T1A, as the top performance yields a 77.1% AUC score. The results on the static prediction are promising and set a good baseline for the real-time prediction scenario, discussed next.

TABLE 4. CONTENT ADOPTION PREDICTION (T2A)

	Precision	Recall	F1	AUC
Logistic Regression	0.433	0.813	0.565	0.755
Random Forests	0.745	0.615	0.674	0.771

4.2.2. T2B: Simulated real-time prediction. We again consider temporal data dependencies to simulate a real-time prediction for T2. Similarly to T1B, in T2B we preserve the temporal ordering of data, and divide the dataset in training and testing according to month-long temporal slices, as summarized by Table 5. *Random Forests* again seems to benefit from the temporal correlations in the data, and the prediction performance at peak improves up to 80.2% AUC. *Logistic Regression* fails again at exploiting temporal information, showing some performance deterioration if compared to T2A. Fig. 2 shows that, for the content adoption prediction, the ranking of the top 11 features in T2B is less stable than that of T1B. The top three most predictive features for this task are (1) ratio of retweets over tweets, (2) number of tweets, and (3) average number of retweets. Note that the latter two top features also appear in the top 3 of the previous task, suggesting an emerging pattern of feature predictive dynamics.

4.3. T3: Predicting interactions with extremists

Our third and last task, namely the prediction of interactions between regular users and extremists, is discussed in the following, again separately for the static (T3A) and real-time (T3B) scenarios.

TABLE 5. REAL-TIME CONTENT ADOPTION PREDICTION (T2B)

Training:	Jan	Jan-Feb	Jan-Mar	Jan-Apr	Jan-May
Testing:	Feb-Jun	Mar-Jun	Apr-Jun	May-Jun	Jun
AUC (LR)	0.682	0.674	0.673	0.703	0.718
Precision (LR)	0.188	0.240	0.148	0.116	0.043
Recall (LR)	0.814	0.367	0.345	0.725	0.362
F1 (LR)	0.305	0.290	0.207	0.199	0.077
AUC (RF)	0.565	0.598	0.676	0.779	0.802
Precision (RF)	0.433	0.384	0.266	0.205	0.070
Recall (RF)	0.087	0.070	0.336	0.648	0.813
F1 (RF)	0.145	0.119	0.297	0.311	0.130

4.3.1. T3A: Time-independent prediction. We expect the interaction prediction task to be the most challenging among the three tasks we proposed. Similarly to content adoption prediction, recent literature has explored the daunting challenge of predicting interaction reciprocity and intensity in social media [20]. Consistently with the prior two tasks, our first approach to interaction prediction is time agnostic: we plan to test whether our system is capable to predict whether a regular user who is mentioned by an extremist account will reply back or not. In this case, positive instances are represented by users who reply to at least one contact initiated by ISIS in the aggregated time period (January-June 2015), whereas negative instances are those regular users who did not reply to any ISIS contact. Table 6 reports the prediction performance of our two models: overall, the task proves challenging as expected, being *Random Forests* the best performer, yielding excellent Recall and 69.2% AUC. *Logistic Regression* provides fair performance with a 65.8% AUC score, and both Precision and Recall around 69%.

TABLE 6. INTERACTION RECIPROCITY PREDICTION (T3A)

	Precision	Recall	F1	AUC
Logistic Regression	0.697	0.690	0.693	0.658
Random Forests	0.686	0.830	0.751	0.692

4.3.2. T3B: Simulated real-time prediction. The final task discussed in this paper is the interaction prediction with temporal data. Given the complexity of this problem, as demonstrated by T3A, we plan to test whether incorporating the temporal dimension will help our models achieve better performance. Table 7 shows that this appears to be the case: *Random Forests* exhibits an improved temporal prediction performance, boasting up to 72.6% AUC, using the first 5 months of data for training, and the last month for

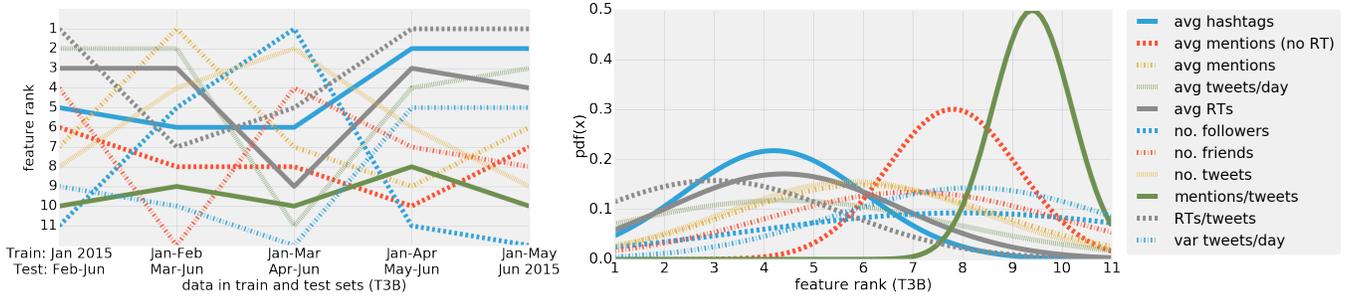


Figure 3. T3B: Feature selection analysis and feature rank distribution (top 11 features)

prediction and evaluation. *Logistic Regression* improves as well, jumping to a 68.3% AUC score. Both models provide very good Precision/Recall performance, if one considers the challenging nature of predicting interaction reciprocity within our context. Fig. 3 summarizes the top 11 features ranking, this time showing a more clear division among top features. The top three features in the interaction reciprocity prediction are (1) ratio of retweets over tweets, (2) average number of hashtags, and (3) average number of retweets. Note that all three features already occurred in the top features of the two previous tasks, reinforcing the notion of a clear pattern of feature predictive power, discussed next.

TABLE 7. REAL-TIME INTERACTION RECIPROcity PREDICTION (T3B)

Training:	Jan	Jan-Feb	Jan-Mar	Jan-Apr	Jan-May
Testing:	Feb-Jun	Mar-Jun	Apr-Jun	May-Jun	Jun
AUC (LR)	0.610	0.589	0.618	0.638	0.683
Precision (LR)	0.562	0.560	0.574	0.553	0.367
Recall (LR)	0.720	0.775	0.813	0.783	0.647
F1 (LR)	0.631	0.650	0.672	0.649	0.468
AUC (RF)	0.628	0.633	0.649	0.671	0.726
Precision (RF)	0.614	0.627	0.603	0.637	0.542
Recall (RF)	0.779	0.676	0.641	0.717	0.765
F1 (RF)	0.687	0.650	0.621	0.675	0.634

4.4. Feature predictive-power analysis

We conclude our analysis by discussing the predictive power of the features adopted by our framework. First, the choice to focus on the top 11 features, rather than the more traditional top 10, is justified by the occurrence of two *ex aequo* in the final ranking of top features, displayed in Table 8. Here, we report the ranking of the top 11 features in the three tasks above. Feature selection is performed on the real-time prediction tasks (not on the time-aggregated ones). This analysis captures the essence of the predictive value of our hand-crafted features in the context of real-time predictions. A clear pattern emerges: (1) ratio of retweets over tweets, (2) average number of hashtags, (2 *ex-aequo*) number of tweets, and (4) average number of retweets, consistently ranked in the top features for the three different prediction tasks. This insight is encouraging: all these features can be easily computed from the metadata reported by the Twitter API, and therefore could be potentially implemented in a real-time detection and prediction system operating on the Twitter stream with unparalleled efficiency.

TABLE 8. FEATURE RANKING ACROSS THE 3 PREDICTION TASKS

Feature	Rank: T1B	T2B	T3B	Final
Ratio of retweets / tweets	4	1	1	1
Avg number of hashtags	2	4	2	2
Number of tweets	1	2	5	=
Avg number of retweets	3	3	3	4
Avg tweets per day	5	8	4	5
Avg no. mentions (w/out retweets)	8	5	8	6
Number of followers	7	6	9	7
Number of friends	6	11	7	8
Avg number of mentions	11	9	6	9
Var tweets / day	9	7	10	=
Ratio of mentions / tweets	10	10	11	11

5. Related Literature

Two relevant research trends recently emerged in the *computational social sciences* and in the *computer science* communities, discussed separately in the following.

Computational social sciences research. This research line is concerned more with understanding the social phenomena revolving around extremist propaganda using online data as a proxy to study individual and group behaviors. Various recent studies focus on English- and Arabic-speaking audiences online to study the effect of ISIS' propaganda and radicalization. One example of the former is the work by Berger and collaborators that provided quantitative evidence of ISIS' use of social media for propaganda. In a 2015 study [21], the authors characterized the Twitter population of ISIS supporters, quantifying its size, provenance, and organization. They argued that most of ISIS' success on Twitter is due to a restricted number of highly-active accounts (500-1000 users). Our analysis illustrates that indeed a limited number of ISIS accounts achieved a very high visibility and followership. Berger's subsequent work [22] however suggested that ISIS' reach (at least among English speakers) has stalled for months as of the beginning of 2016, due to more aggressive account suspension policies enacted by Twitter. Again, a limited amount of English accounts sympathetic to ISIS was found (less than one thousand), and these users were mostly interacting with each other, while being only marginally successful at acquiring other users' attention. This analysis suggests a mechanism of diminishing returns for extremist social media propaganda.

Using Twitter data as a historical archive, some researchers [23] recently tried to unveil the roots of support for

ISIS among the Arabic-speaking population. Their analysis seems to suggest that supporters of the extremist group have been discussing about Arab Spring uprisings in the past significantly more than those who oppose ISIS on Twitter. Although their method to separate ISIS supporters from opposers is simplistic, the findings relating narrative framing and recruitment mechanisms are compatible with the literature on social protest phenomena [24], [25], [26].

A few studies explored alternative data sources: one interesting example is the work by Vergani and Bliuc [27] that uses sentiment analysis (Linguistic Inquiry and Word Count [28]) to investigate how language evolved across the first 11 Issues of Dabiq, the flagship ISIS propaganda magazine. Their analysis offers some insights about ISIS radicalization motives, emotions and concerns. For example, the authors found that ISIS has become increasingly concerned with females, reflecting their need to attract women to create their utopia society, not revolving around warriors but around families. ISIS also seems to have increased the use of internet jargon, possibly to connect with the identities of young individuals online.

Computer science research. This research stream concerns more with the machine learning and data aspects, to model, detect, and/or predict social phenomena such as extremism or radicalization often with newly-developed techniques.

One of the first computational frameworks, proposed by Bermingham *et al.* [29] in 2009, combined social network analysis with sentiment detection tools to study the agenda of a radical YouTube group: the authors examined the topics discussed within the group and their polarity, to model individuals' behavior and spot signs of extremism and intolerance, seemingly more prominent among female users. The detection of extremist content (on the Web) was also the focus of a 2010 work by Qi *et al.* [30]. The authors applied hierarchical clustering to extremist Web pages to divide them into different categories (religious, politics, etc.).

Scanlon and Gerber proposed the first method to detect cyber-recruitment efforts in 2014 [31]. They exploited data retrieved from the Dark Web Portal Project [32], a repository of posts compiled from 28 different online fora on extremist religious discussions (e.g., Jihadist) translated from Arabic to English. After annotating a sample of posts as recruitment efforts or not, the authors use Bayesian criteria and a set of textual features to classify the rest of the corpus, obtaining good accuracy, and highlighted the most predictive terms.

Along the same trend, Agarwal and Sureka proposed different machine learning strategies [33], [34], [35], [36] aimed at detecting radicalization efforts, cyber recruitment, hate promotion, and extremist support in a variety of online platforms, including YouTube, Twitter and Tumblr. Their frameworks leverage features of contents and metadata, and combinations of crawling and unsupervised clustering methods, to study the online activity of Jihadist groups on the platforms mentioned above.

Concluding, two very recent articles [37], [38] explore the activity of ISIS on social media. The former [37] focuses on Twitter and aims at detecting users who exhibit signals of behavioral change in line with radicalization: the

authors suggest that out of 154K users only about 700 show significant signs of possible radicalization, and that may be due to social homophily rather than the mere exposure to propaganda content. The latter study [38] explores a set of 196 pro-ISIS aggregates operating on VKontakte (the most popular Russian online social network) and involving about 100K users, to study the dynamics of survival of such groups online: the authors suggest that the development of large and potentially influential pro-ISIS groups can be hindered by targeting and shutting down smaller ones. For additional pointers we refer the interested reader to two recent literature reviews on this topic [39], [40].

6. Conclusions

In this article we presented the problem of predicting online extremism in social media. We defined three machine learning tasks, namely the detection of extremist users, the prediction of extremist content adoption, and the forecasting of interactions between extremist users and regular users. We tapped into the power of a crowd-sourcing project that aimed at manually identifying and reporting suspicious or abusive activity related to ISIS radicalization and propaganda agenda, and collected annotations to build a ground-truth of over 25 thousand suspended Twitter accounts. We extracted over three million tweets related to the activity of these accounts in the period of time between January and June 2015. We also randomly identified an equal-sized set of regular users exposed to the extremist content generated by the ISIS accounts, and collected almost 30 million tweets generated by the regular users in the same period of time.

By means of state-of-the-art learning models we managed to accomplish predictions in two types of scenarios, a static one that ignores temporal dependencies, and a simulated real-time case in which data are processed for training and testing by respecting the timeline of content availability. The two learning models, and the set of 52 features that we carefully crafted, proved very effective in all of the six combinations of forecasts (three prediction tasks each with two prediction modalities, static and real-time). The best performance in terms of AUC ranges between 72% and 93%, depending on the complexity of the considered task and the amount of training data available to the system.

We concluded our analysis by investigating the predictive power of different features. We focused on the top 11 most significant features, and we discovered that some of them, such as the ratio of retweets to tweets, the average number of hashtags adopted, the sheer number of tweets, and the average number of retweets generated by each user, systematically rank very high in terms of predictive power. Our insights shed light on the dynamics of extremist content production as well as some of the network and timing patterns that emerge in this type of online conversation.

Our work is far from concluded: for the future, we plan to identify more realistic and complex prediction tasks, to analyze the network and temporal dynamics of extremist discussion, and to deploy a prototype system that allows for real-time detection of signatures of abuse on social media.

Acknowledgments

This work has been supported in part by the Office of Naval Research (grant no. N15A-020-0053). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

References

- [1] Q. Schiermeier, "Terrorism: Terror prediction hits limits." *Nature*, vol. 517, no. 7535, p. 419, 2015.
- [2] S. Reardon, "Terrorism: science seeks roots of terror," *Nature*, vol. 517, no. 7535, pp. 420–421, 2015.
- [3] J. Berger and B. Strathearn, "Who matters online: measuring influence, evaluating content and countering violent extremism in online social networks," *Int. Centre for the Study of Radicalisation*, 2013.
- [4] A. Fisher, "How jihadist networks maintain a persistent online presence," *Perspectives on Terrorism*, vol. 9, no. 3, 2015.
- [5] J. Stern and J. M. Berger, *ISIS: The state of terror*. Harper, 2015.
- [6] P. Cockburn, *The rise of Islamic State: ISIS and the new Sunni revolution*. Verso Books, 2015.
- [7] M. Weiss and H. Hassan, *ISIS: Inside the army of terror*. Simon and Schuster, 2015.
- [8] C. A. Davis, O. Varol, E. Ferrara, A. Flammini, and F. Menczer, "Botornot: A system to evaluate social bots," in *Proceedings of the 25th International Conference Companion on World Wide Web*, 2016.
- [9] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, "Scikit-learn: Machine learning in python," *The Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [10] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [11] T. Hastie, R. Tibshirani, J. Friedman, and J. Franklin, "The elements of statistical learning: data mining, inference and prediction," *The Mathematical Intelligencer*, vol. 27, no. 2, pp. 83–85, 2005.
- [12] A. Mislove, S. Lehmann, Y.-Y. Ahn, J.-P. Onnela, and J. N. Rosenquist, "Understanding the demographics of Twitter users," in *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media*, 2011.
- [13] E. Ferrara, O. Varol, C. Davis, F. Menczer, and A. Flammini, "The rise of social bots," *arXiv preprint arXiv:1407.5225*, 2014.
- [14] R. Ghosh, T. Surachawala, and K. Lerman, "Entropy-based classification of retweeting activity on twitter," in *Proceedings of KDD workshop on Social Network Analysis (SNA-KDD)*, August 2011.
- [15] V. Subrahmanian, A. Azaria, S. Durst, V. Kagan, A. Galstyan, K. Lerman, L. Zhu, E. Ferrara, A. Flammini, F. Menczer *et al.*, "The DARPA Twitter bot challenge," *arXiv preprint arXiv:1601.05140*, 2016.
- [16] J. Ratkiewicz, M. Conover, M. Meiss, B. Goncalves, A. Flammini, and F. Menczer, "Detecting and tracking political abuse in social media," in *Proceedings of the 5th International AAAI Conference on Weblogs and Social Media*, 2011, pp. 297–304.
- [17] E. Ferrara, O. Varol, F. Menczer, and A. Flammini, "Detection of promoted social media campaigns," in *Proceedings of the 10th International Conference on Web and Social Media*, 2016.
- [18] I. Guyon and A. Elisseeff, "An introduction to variable and feature selection," *The Journal of Machine Learning Research*, vol. 3, pp. 1157–1182, 2003.
- [19] K. Lerman and R. Ghosh, "Information Contagion: An Empirical Study of the Spread of News on Digg and Twitter Social Networks," in *Proceedings of the 4th International AAAI Conference on Weblogs and Social Media*, 2010, pp. 90–97.
- [20] E. Gilbert and K. Karahalios, "Predicting tie strength with social media," in *Proceedings of the SIGCHI conference on human factors in computing systems*. ACM, 2009, pp. 211–220.
- [21] J. Berger and J. Morgan, "The ISIS twitter census: Defining and describing the population of isis supporters on twitter," *The Brookings Project on US Relations with the Islamic World*, vol. 3, no. 20, 2015.
- [22] J. Berger and H. Perez, "The Islamic States diminishing returns on Twitter," *GW Program on extremism*, vol. 02-16, 2016.
- [23] W. Magdy, K. Darwish, and I. Weber, "#failedrevolutions: Using Twitter to study the antecedents of ISIS support," *First Monday*, vol. 21, no. 2, 2016.
- [24] S. González-Bailón, J. Borge-Holthoefer, A. Rivero, and Y. Moreno, "The dynamics of protest recruitment through an online network," *Scientific reports*, vol. 1, 2011.
- [25] M. D. Conover, C. Davis, E. Ferrara, K. McKelvey, F. Menczer, and A. Flammini, "The geospatial characteristics of a social movement communication network," *PloS one*, vol. 8, no. 3, e55957, 2013.
- [26] M. D. Conover, E. Ferrara, F. Menczer, and A. Flammini, "The digital evolution of occupy wall street," *PloS one*, vol. 8, no. 5, 2013.
- [27] M. Vergani and A.-M. Bliuc, "The evolution of the ISIS' language: a quantitative analysis of the language of the first year of dabiq magazine," *Security, Terrorism, and Society*, vol. 1, no. 2, 2015.
- [28] Y. R. Tausczik and J. W. Pennebaker, "The psychological meaning of words: LIWC and computerized text analysis methods," *Journal of language and social psychology*, vol. 29, no. 1, pp. 24–54, 2010.
- [29] A. Bermingham, M. Conway, L. McInerney, N. O'Hare, and A. F. Smeaton, "Combining social network analysis and sentiment analysis to explore the potential for online radicalisation," in *2009 International Conference on Advances in Social Network Analysis and Mining (ASONAM)*. IEEE, 2009, pp. 231–236.
- [30] X. Qi, K. Christensen, R. Duval, E. Fuller, A. Spahiu, Q. Wu, and C.-Q. Zhang, "A hierarchical algorithm for clustering extremist web pages," in *2010 International Conference on Advances in Social Networks Analysis and Mining (ASONAM)*, 2010, pp. 458–463.
- [31] J. R. Scanlon and M. S. Gerber, "Automatic detection of cyber-recruitment by violent extremists," *Security Informatics*, vol. 3, no. 1, pp. 1–10, 2014.
- [32] H. Chen, W. Chung, J. Qin, E. Reid, M. Sageman, and G. Weimann, "Uncovering the dark web: A case study of jihad on the web," *Journal of the American Society for Information Science and Technology*, vol. 59, no. 8, pp. 1347–1359, 2008.
- [33] S. Agarwal and A. Sureka, "A focused crawler for mining hate and extremism promoting videos on youtube," in *Proceedings of the 25th ACM conference on Hypertext and social media*, 2014, pp. 294–296.
- [34] A. Sureka and S. Agarwal, "Learning to classify hate and extremism promoting tweets," in *Intelligence and Security Informatics Conference (JISIC), 2014 IEEE Joint*. IEEE, 2014, pp. 320–320.
- [35] S. Agarwal and A. Sureka, "Using knn and svm based one-class classifier for detecting online radicalization on twitter," in *Distributed Computing and Internet Technology*. Springer, 2015, pp. 431–442.
- [36] —, "Spider and the flies: Focused crawling on tumblr to detect hate promoting communities," *arXiv preprint arXiv:1603.09164*, 2016.
- [37] M. Rowe and H. Saif, "Mining pro-ISIS radicalisation signals from social media users," in *Proceedings of the 10th International Conference on Web and Social Media*, 2016.
- [38] N. Johnson, M. Zheng, Y. Vorobyeva, A. Gabriel, H. Qi, N. Velasquez, P. Manrique, D. Johnson, E. Restrepo, C. Song *et al.*, "New online ecology of adversarial aggregates: Isis and beyond," *arXiv preprint arXiv:1603.09426*, 2016.
- [39] D. Correa and A. Sureka, "Solutions to detect and analyze online radicalization: a survey," *arXiv preprint arXiv:1301.4916*, 2013.
- [40] S. Agarwal and A. Sureka, "Applying social media intelligence for predicting and identifying on-line radicalization and civil unrest oriented threats," *arXiv preprint arXiv:1511.06858*, 2015.