

Predicting enterprise cyber incidents using social network analysis on dark web hacker forums

Soumajyoti Sarkar

Arizona State University
Tempe, Arizona, U.S.A.

Mohammad Almukaynizi

Arizona State University
Tempe, Arizona, U.S.A.

Jana Shakarian

Cyber Reconnaissance, Inc.
Tempe, Arizona, U.S.A.

Paulo Shakarian

Arizona State University
Tempe, Arizona, U.S.A.

ABSTRACT

With the rise in security breaches over the past few years, there has been an increasing need to mine insights from social media platforms to raise alerts of possible attacks in an attempt to defend conflict during competition. We use information from dark web forums by leveraging the reply network structure of user interactions with the goal of predicting enterprise cyberattacks. We use a suite of social network features on top of supervised learning models and validate them using a binary classification problem that attempts to predict whether there would be an attack on any given day for an organization. We conclude from our experiments, which gathered information from 53 forums on the dark web over a span of 12 months and attempted to predict real-world cyberattacks across 2 security incidents, that analyzing the path structure between groups of users is better than merely studying centralities like Pagerank or relying on user-posting statistics in forums.

INTRODUCTION

With recent data breaches at organizations such as Yahoo, Uber, and Equifax¹ emphasizing the increasing financial and social impacts of cyberattacks, there has been an enormous requirement for technologies that could alert such organizations to possible data breaches. These breaches are a direct or indirect result of cyber, electronic, and information operations to infiltrate systems and infrastructure as well as gain unauthorized access to information, thus setting an example of conflict during competition. On the vulnerability

© 2019 Soumajyoti Sarkar, Mohammad Almukaynizi, Jana Shakarian, Paulo Shakarian

Some of the authors are supported through the AFOSR Young Investigator Program (YIP) grant FA9550-15-1-0159, ARO grant 11NF-15-1-0282, and the DoD Minerva program grant N00014-16-1-2015.

¹ <https://www.consumer.ftc.gov/blog/2017/09/equifax-data-breach-what-do> <https://www.consumer.ftc.gov/blog/2016/09/yahoo-breach-watch>

front, Risk Based Security's VulnDB database² published a total of 4,837 vulnerabilities in a quarter of 2017, which was around 30 percent higher than the previous year. This motivates the need for extensive systems that can utilize vulnerability-associated information from external sources to alert organizations to such cyberattacks. The dark web is one such place on the internet where users can share information on software vulnerabilities and ways to exploit them^{[1],[15]}. Surprisingly, it might be difficult to track the actual intentions of those users, thus making it necessary to use data mining and learn to identify the discussions among the noise that could potentially raise alerts on attacks on external enterprises. In this paper, we leverage the information obtained from analyzing the reply network structure of discussions in dark web forums to understand the extent to which dark web information can be useful for predicting real-world cyberattacks.

Most of the work on vulnerability discussions on trading exploits in underground forums^{[9],[13],[14]} and related social media platforms like Twitter^{[2],[8],[15]} have focused on two aspects: (1) analyzing vulnerabilities discussed or traded in the forums and the markets, thereby giving rise to the belief that the "life cycle of vulnerabilities" in these forums and marketplaces and their exploitation have a significant impact on real-world cyberattacks^{[13],[14]}; and (2) prioritizing or scoring vulnerabilities using these social media platforms or binary file appearance logs of machines to predict the risk state of machines or systems^{[7],[11]}. These two components have been used in silos; however, and in this paper, we ignore the steps between vulnerability exploit analysis and the final task of real-world cyberattack prediction by removing the preconceived notions used in earlier studies where vulnerability exploitation was considered a precursor towards attack prediction. We instead hypothesize about user interaction dynamics conceived through posts surrounding these vulnerabilities on these underground platforms to generate warnings for future attacks. We note that we *do not* consider whether vulnerabilities have been exploited in these discussions since a lot of zero-day attacks^[11] might occur before such vulnerabilities are even indexed and their gravity might lie hidden in discussions related to other associated vulnerabilities or some discussion on exploits. We based our research on the dynamics of all kinds of discussions on dark web forums; however, we attempted to filter out the noise to mine important patterns by examining whether pieces of information gained traction within important communities.

To this end, the major contributions of this research investigation are as follows:

- ◆ We create a network mining technique using the directed reply network of users who participate in dark web forums to extract a set of specialized users we term *experts* whose posts with *popular vulnerability mentions* gain attention from other users in a specific time frame.
- ◆ Following this, we generate several time series of features that capture the dynamics of interactions centered around these *experts* across individual forums as well as general feature time series based on social network and forum posting statistics.

² <https://www.riskbasedsecurity.com/2017/05/29-increase-in-vulnerabilitiesalready-disclosed-in-2017/>

- ◆ We use these time series features and train a supervised learning model based on logistic regression with attack labels for two different incidents from an organization to predict daily attacks. We obtain the best results with an F1 score of 0.53 on a feature that explores the path structure between *experts* and other users compared to the random (without prior probabilities) F1 score of 0.37. Additionally, we identify instances of superior feature performance based on discussions involving vulnerability information rather than network centralities and forum posting statistics.

The rest of the paper is organized as follows: We introduce several terms and the dataset related to the vulnerabilities and dark web in section II; the general framework for attack prediction, including feature curation and learning models, in section III; and, finally, our experimental evaluations in section IV.

II. BACKGROUND AND DATASET

In this section, we describe the dataset that we used in our research to analyze the interaction of patterns of dark web users and the real-world security incident³ data that we used as ground truth (GT) for the evaluation of our prediction models.

A. Enterprise-Relevant External Threats (Ground Truth (GT))

Our GT was based on data on cyberattacks on Armstrong Corporation systems occurring between April 2016 and September 2017; we obtained this data from the Intelligence Advanced Research Projects Activity Cyberattack Automated Unconventional Sensor Environment program⁴. Some of the relevant attributes in this data are “event type” and “event occurred date.” Event type is the type of attack and event occurred date is the date on which a particular attack occurred. The event types examined in this study are “malicious email” and “endpoint malware.” Malicious email refers to an incident in which an individual in the organization received an email that contained a malicious attachment or link. Endpoint malware refers to malware discovered on an endpoint device. This includes, but is not limited to, ransomware, spyware, and adware. As shown in figure 1, the distribution of attacks over time is different for the events. The total number of incidents reported for the events are as follows: 119 tagged as *endpoint-malware* and 135 as malicious-email events, resulting in a total of 280 incidents over a span of 17 months that were considered in our study.

B. Dark web data

The dark web forms a small part of the deep web, the part of the web not indexed by web search engines, although sometimes the term “deep web” is mistakenly used to refer to the dark web. We obtained all dark web data used in this study through an application programming interface provided by a commercial platform⁵.

³ We often use the terms “attacks,” “incidents,” and “events” interchangeably.

⁴ <https://www.iarpa.gov/index.php/research-programs/cause>

⁵ Data is provided by Cyber Reconnaissance, Inc., www.cyr3con.ai

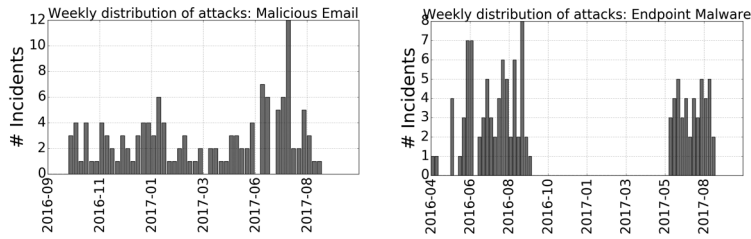


Fig. 1: Weekly occurrence of security breach incidents of different types (a) malicious email and (b) endpoint malware

The structure of a dark web forum is hierarchical: each forum consists of several independent threads, a thread caters to a particular discussion on a topic, and each thread spans several posts initiated by multiple users over time. We note that one user can appear multiple times in the sequence of posts depending on when and how many times the user posted in that thread. However, the dataset we obtained does not contain the hierarchical information of reposting - it does not provide us with which user a particular user replied to, while posting or replying in a thread. We filter out forums based on a threshold number of posts that were created in the time frame of January 2016 to September 2017. We gathered data from 179 forums in which the total number of unique posts was 557,689, irrespective of the thread to which they belonged. The number of forums with less than 100 posts was large and, therefore, we only considered forums that had greater than 5,000 posts during that time period, which gave us a total of 53 forums. We denote the set of 53 forums used in this dataset with the symbol F .

Common Vulnerabilities and Exposures (CVEs): The database of CVEs maintained on a platform operated by the MITRE Corporation⁶ provides identity mapping for publicly known information security vulnerabilities and exposures. We collect all of the information regarding the vulnerability mentions in dark web forums during the period between January 2016 and October 2017. The total number of CVEs mentioned in the posts across all forums during this period was 3,553.

CVE - Common Platform Enumeration (CPE) Mapping: A CPE is a structured naming scheme for identifying and grouping clusters of information technology systems, software, and packages maintained on the National Vulnerability Database (NVD) platform operated by the National Institute of Standards and Technology (NIST)⁷. Each CVE can be assigned to different CPE groups based on the naming system of CPE families, as described in [9]. Similarly, each CPE family can have several CVEs that conform to its vendors and products to which the CPE caters. In order to cluster the set of CVEs in our study into a set of CPE groups, we use the set of CPE tags for each CVE from the NVD database maintained by NIST. For the CPE tags, we only consider the operating system platform and the application environment tags for each unique CPE. Examples of CPEs include Microsoft Windows_95, Canonical ubuntu_linux, and, Hp elitebook_725_g3. The first component in each of these CPEs denotes the operating system platform and the second component denotes the application environment and its versions.

⁶ <http://cve.mitre.org>
⁷ <https://nvd.nist.gov/cpe.cfm>

III. FRAMEWORK FOR ATTACK PREDICTION

The mechanism for attack predictions can be described in three steps: (1) Given a time point t for which we need to predict an enterprise attack of a particular event type, (2) we use features from the dark web forums prior to t and (3) we use these features as input for a learned model that predicts attacks on t . So one of the main tasks involves learning the attack prediction model for each event type. Below we describe steps (2) and (3) - curating features and building supervised learning models.

A. Curating Features

We first describe the mechanism in which we build temporal networks, following which we describe the features used for the prediction problem. We build three groups of features across forums: (1) Expert-centric; (2) User/forum statistics; and (3) Network centralities.

Dark Web Reply Network: We assume the absence of global user identification (IDs) across forums⁸ and therefore analyze the social interactions using networks induced on specific forums instead of considering the global network across all forums. We denote the directed reply graph of a forum $f \in F$ by $G^f = (V^f, E^f)$ where V^f denotes the set of users who posted or replied in some thread in forum f at some time in our considered time frame of data; and E^f denotes the set of three-tuple (u_1, u_2, rt) directed edges where $u_1, u_2 \in V^f$ and rt denotes the time at which u_1 replied to a post of u_2 in some thread in f , with $u_1 \rightarrow u_2$ denoting the edge direction. We denote by $G_\tau^f = (V_\tau^f; E_\tau^f)$, a temporal subgraph of G^f , τ being a time window such that V_τ^f denotes the set of individuals who posted in f in that window; in addition, E_τ^f denotes the set of tuples $(v_1; v_2; rt)$ such that $rt \in \tau, v_1; v_2 \in V_\tau^f$. We use two operations to create temporal networks: "Create," which takes a set of forum posts in f within a time window τ as input and creates a temporal subgraph G_τ^f and "Merge," which takes two temporal graphs as input and merges them to form an auxiliary graph. To keep the notations simple, we drop the symbol f when we describe the operations for a specific forum in F as context but which does apply for any forum $f \in F$. We describe these two operations in brief; however, a detailed algorithm relating the network construction is given in algorithm 1 of appendix A1.⁹ We adopt an incremental analysis approach by splitting the entire set of time points in our frame of study into a sequence of time windows $\Gamma = \{\tau_1, \tau_2, \dots, \tau_Q\}$, where each subsequence $\tau_i, i \in [1, Q]$ is equal in time span and non-overlapping and the subsequences are ordered by their starting time points for their respective span.

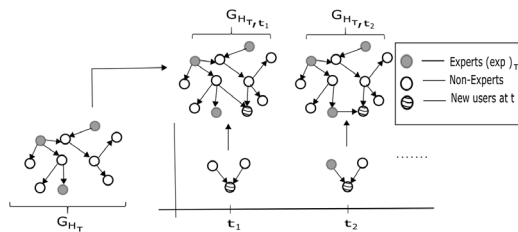


Fig. 2: An illustration to show the Merge operation: $G_{H, \tau}$ denotes the historical network which is used to compute the experts shown in gray. $\{G_{t_1}, G_{t_2}, \dots\}$ denote the networks at time $t_1, t_2, \dots \in \tau, t \in \Gamma$.

⁸ Note that even in the presence of global user IDs across forums, a lot of anonymous or malicious users would create multiple profiles across forums and create multiple posts with profiles; identifying and merging such profiles is an active area of research.

⁹ Online appendix: <http://www.public.asu.edu/~ssarka18/appendix.pdf>

CREATE: *Creating the reply graph* - Let h be a particular thread or topic within a forum f containing posts by users $V_h^f = \{u_1, \dots, u_k\}$ posted at corresponding times $T_h^f = \{t_1, \dots, t_k\}$, where k denotes the number of posts in that thread and $t_i < t_j$ for any $i > j$; that is, the posts are chronologically ordered. To create the set of edges E_f , we connect two users $(u_i, u_j) \in V_h^f$ such that $i > j$; that is, user u_i has potentially replied to u_j , and is subject to a set of spatial and temporal constraints (appendix A1). These constraints make up for the absence of exact information about the reply as to whom u replied to in a particular post in h .

MERGE: *Merging network* - In order to create a time series feature $T_{x,f}$ for feature x from threads in forum f that maps each time point $t \in \tau$, $\tau \in \Gamma$ to a real number, we use two networks: (1) the historical network G_{H_τ} , which spans over time H_τ such that $\forall t \in H_\tau$, and $t \in \tau$, so that we have $t' < t$; and (2) the network G_t^f induced by user interactions between users in E_t , which varies temporally for each $t \in \tau$. We note that the historical network G_{H_τ} is different for each subsequence and same for all $t \in \tau$, so that as the subsequences $\tau \in \Gamma$ progress with time, the historical network G_{H_τ} also changes; in addition, we discuss the choice of spans $\tau \in \Gamma$ and H_τ in section IV. Finally, for computing feature values for each time point $t \in \tau$, we merge the two networks G_{H_τ} and G_t^f to form the auxiliary network $G_{H_\tau,t} = (V_{H_\tau,t}, E_{H_\tau,t})$, where $V_{H_\tau,t} = V_{H_\tau} \cup V_t$ and $E_{H_\tau,t} = E_{H_\tau} \cup E_\tau$. A visual illustration of this method is shown in figure 2. Now we describe the several features we used that would be fed to a learning model for attack prediction. We compute time series of several features x , $T_{x,f}[t]$ for every time point t in our frame of study and for every forum f separately.

1. Expert-Centric Features

We extract a set of users we term *experts* who have a history of CVE mentions in their posts and whose posts have gained attention in terms of replies. Following that, we mine several features that explain how attention is broadcast by these experts to other posts. All of these features are computed using the auxiliary networks $G_{H_\tau,t}$ for each time t . Our hypothesis is based on the premise that any unusual activity must spur attention from users who have knowledge about vulnerabilities.

We focus on users whose posts in a forum contain the most-discussed CVEs belonging to important CPEs during the time frame of analysis, where the importance will shortly be formalized. For each forum f , we use the historical network $G_{H_\tau}^f$ to extract the set of *experts* relevant to time frame τ ; that is, $exp_t^f \in V_{H_\tau}^f$. First, we extract the top CPE groups CPTop in the time frame H based on the number of historical mentions of CVEs. We sort the CPE groups based on the sum of the CVE mentions that belong to the respective CPE groups and take the top five CPE groups by the sum in each H_τ . Using these notations, the experts exp_t^f from history H_τ considered for time span τ are defined as users in f with the following three constraints:

- (1) Users who have mentioned a CVE in their posts in H_τ . This ensures that the user engages in the forums with content that is relevant to vulnerabilities.

- (2) Let $\theta(u)$ denote the set of CPE tags of the CVEs mentioned by user u in his/her posts in H_τ , such that it follows either $\theta(u) \in CP_\tau^{top}$, where the user’s CVEs are grouped in less than five CPEs, or, $CP_\tau^{top} \in \theta(u)$ in cases where a user has posts with CVEs in the span H_τ grouped into more than five CPEs. This constraint filters out users who discuss vulnerabilities that are not among the top CPE groups in H_τ .
- (3) The in-degree of the user u in GH_τ should cross a threshold. This constraint ensures that there are a significant number of users who potentially responded to this user, thus establishing u ’s central position in the reply network. Essentially, the set of experts exp from H_τ would be used for all of the time points in τ .

We curate path- and community-based features based on the experts listed in table I. These expert-centric features try to quantify the distance between an expert and a daily user (non-expert) in terms of how fast a post from that user receives attention from the expert. In that sense, the community features also measure the like-mindedness of non-experts and experts.

Group	Features	Description
Expert centric	Graph Conductance	$\tau_x[t] = \frac{1}{\sum_{e \in experts} \sum_{u \in V_t \setminus experts} \pi(expert) P_{xy}}$ where $\pi(\cdot)$ is the stationary distribution of the network $G_{H_\tau, t}$. P_{xy} denotes the probability of random walk from vertices x to y . The conductance represents the probability of taking a random walk from any of the <i>experts</i> to one of the users in $V_t \setminus experts$, normalized by the probability weight of being on an expert.
	Shortest Path	$\tau_x[t] = \frac{1}{ experts } \sum_{e \in experts} \min_{u \in V_t \setminus experts} s_{e,u}$ where $s_{e,u}$ denotes the shortest path from an expert e to user u following the direction of edges.
	Expert replies	$\tau_x[t] = \frac{1}{ experts } \sum_{e \in experts} OutNeighbors(e) $ where $OutNeighbors(\cdot)$ denotes the out neighbors of user in the network $G_{H_\tau, t}$.
	Common Communities	$\tau_x[t] = \mathcal{N}(c(u)) \mid c(u) \in c_{experts} \wedge u \in V_t \setminus experts $ where $c(u)$ denotes the community index of user u , $c_{experts}$ that of the experts and $\mathcal{N}(\cdot)$ denotes a counting function. It counts the number of users who share communities with experts.
Forum/User Statistics	Number of threads	$\tau_x[t] = \{h \mid \text{thread } h \text{ was posted on } t\} $
	Number of users	$\tau_x[t] = \{u \mid \text{user } u \text{ posted on } t\} $
	Number of expert threads	$\tau_x[t] = \{h \mid \text{thread } h \text{ was posted on } t \text{ by users } u \in \text{experts}\} $
	Number of CVE mentions	$\tau_x[t] = \{CVE \mid \text{CVE was mentioned in some post on } t\} $
Network Centralities	<i>Outdegree</i> _k	$\tau_x[t] = \text{Average value of top } k \text{ users, by outdegree on } t$
	<i>Outdegree</i> _k CVE	$\tau_x[t] = \text{Average value of top } k \text{ users with more than } 1 \text{ CVE mention in their posts, by outdegree on } t$
	<i>Pagerank</i> _k	$\tau_x[t] = \text{Average value of top } k \text{ users, by Pagerank on } t$
	<i>Pagerank</i> _k CVE	$\tau_x[t] = \text{Average value of top } k \text{ users with more than } 1 \text{ CVE mention in their posts, by pagerank on } t$
	<i>Betweenness</i> _k	$\tau_x[t] = \text{Average value of top } k \text{ users, by Betweenness on } t$
<i>Betweenness</i> _k CVE	$\tau_x[t] = \text{Average value of top } k \text{ users with more than } 1 \text{ CVE mention in their posts, by betweenness on } t$	

Fig. 1: Weekly occurrence of security breach incidents of different types (a) malicious email and (b) endpoint malware

Why focus on experts? To show the significance of these properties in comparison those of other users, we examine the time periods of 3 widely known security events: the Wannacry ransomware attack that happened on May 12, 2017, and the vulnerability MS-17-010; the Petya cyberattack on June 27, 2017, with the associated vulnerabilities CVE-2017-0144, CVE-2017-0145, and MS-17-010; and the Equifax breach attack that occurred primarily on March 9, 2017, and vulnerability CVE-2017-5638. We consider two sets of users across all forums - exp_τ , where G_{H_τ} denotes the corresponding historical network prior to τ in which these three events occurred and the second set of users being all U_{alt} who are not experts and who fail either one of two constraints: they have mentioned CVEs in their posts which do not belong to CP^{top} or their in-degree in G_{H_τ} lies below the threshold. We consider G_{H_τ} being induced by users in the last 3 weeks prior to the occurrence week of each event for both cases and we consider the

total number of interactions, ignoring the direction of reply of these users with other users. Let deg_{exp} denote the vector of counts of interactions in which the experts were involved and deg_{alt} denote the vector of counts of interactions in which the users in U_{alt} were involved. We randomly pick a number of users from U_{alt} equal to the number of experts and sort the vectors by count. We conduct a two-sample t -test on the vectors deg_{exp} and deg_{alt} . The null hypothesis H_0 and the alternate hypothesis H_1 are defined as follows:

- ◆ $H_0 : \text{deg}_{\text{exp}} \leq \text{deg}_{\text{alt}}$
- ◆ $H_1 : \text{deg}_{\text{exp}} > \text{deg}_{\text{alt}}$

The null hypothesis is rejected at significance level $\alpha = 0.01$ with a p-value of 0.0007. This suggests that with high probability, experts tend to interact more prior to important, real-world cybersecurity breaches than other users who randomly post CVEs.

Now, we conduct a second t -test where we randomly pick 4 weeks not in the weeks considered for the data breaches to pick users U_{alt} with the same constraints. We use the same hypotheses as above and, when we perform statistical tests for significance, we find that the null hypothesis is not rejected at $\alpha=0.01$ with a p -value close to 0.05. This empirical evidence from the t -test also suggests that the interactions with exp are more correlated with an important cybersecurity incident than those of users who post CVEs not in top CPE groups and that, therefore, it is better to focus on users exhibiting our desired properties as experts for cyberattack prediction. Note that the t -test evidence also incorporates a special temporal association since we collected events from three interleaved time frames corresponding to the event dates.

2. User/Forum Statistics Features

We try to see whether the forum or user-posting statistics are themselves indicators of future cyberattacks; for this, we compute Forum/User Statistics, as described in table I.

3. Network Centrality Features

In addition, we also tested several network centrality features mentioned in table I. The purpose is to check whether the emergence of central users in the reply network $G_t, t \in \tau$, is a good predictor of a cyberattack. We note that, in this case, we only use the daily reply networks to compute the features, unlike in the case of the expert-centric network features, where we use $G_{H\tau, t}$.

B. Building Supervised Learning Models

In this section we explain how we use the time series data $T_{x,f}$ to predict an attack at any given time point t . We consider a supervised learning model in which the time series T_x is formed by averaging $T_{x,f}$ across all forums in $f \in F$ at each time point t and then used for the prediction task. We treat the attack prediction problem in this paper as a binary classification problem in which the objective is to predict whether there would be an attack at a given time point t . Since the incident data in this paper contains the number of incidents that

occurred at time point t , we assign a label of 1 for t if there was at least one attack at t and 0 otherwise.

In [4], the authors studied the effect of longitudinal sparsity in high-dimensional time series data. They propose assigning weights to the same features at different time spans to capture the temporal redundancy. We use two parameters: β , which denotes the start time prior to t from where we consider the features for prediction, and n , the time span for the features to be considered. In figure 3, to predict an attack occurrence at time t , we use the features for each time $t_h \in [t_{-n-\beta}, t_{-\beta}]$. Here we use logistic regression with longitudinal ridge sparsity that models the probability of an attack as follows with X being the set of features and β being the vector of coefficients:

$$P(attack(t) = 1 | \mathbf{X}) = \frac{1}{1 + e^{-(\beta_0 + \sum_{k=\eta+\delta}^{\delta} \beta_k x_{t-k})}} \tag{1}$$

The final objective function to minimize over N instances where N is the number of time points spanning the attack time frame: $l(\beta) = -\sum_{i=1}^N (y_i(\beta_0 + \mathbf{x}_i^T \beta) - \log(1 + \exp^{\beta_0 + \mathbf{x}_i^T \beta}) + \lambda \beta^T \beta$, with y being the instance label.

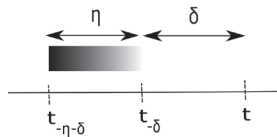


Fig. 3: Temporal feature selection window for predicting an attack at time t

One of the major problems of the dataset is the imbalance in the training and test dataset, as will be described in section IV; thus, in order to use all features in each group together for prediction, we use three additional regularization terms: the L1 penalty, the L2 penalty, and the group lasso regularization [5]. The final objective function can be written as:

$$l(\beta) = -\sum_{i=1}^N \log(1 + e^{-y_i(\beta^T \mathbf{x}_i)}) + \frac{m}{2} \|\beta\|_2^2 + l \|\beta\|_1 + g \cdot GL(\beta) \tag{2}$$

where m , l , and g are the hyper-parameters for the regularization terms and the $GL(\beta)$ term is $\sum_{g=1}^G \|\beta_{I_g}\|_2$, where I_g is the index set belonging to the g^{th} group of variables, $g = 1 \dots G$. Here each g is the time index $t_h \in [t_{-n-\beta}, t_{-\beta}]$, so this group variable selection selects all features of one time in history while reducing some other time points to 0. It has the attractive property that it performs variable selection at the temporal group level and is invariant under (group-wise) orthogonal transformations, like ridge regression. We note that while there are several other models that could be used for prediction that incorporate the temporal and sequential nature of the data, like hidden markov models and recurrent neural networks, the logit model allows us to transparently adjust to the sparsity of data, especially in the absence of a large dataset.

IV. EXPERIMENTAL EVALUATIONS

In our work, the granularity for each time index in the T function is 1 day; that is, we compute feature values over all days in the time frame of our study. For incrementally computing the values of the time series, we consider the time span of each subsequence $\tau \in T$ as 1 month, and for each τ , we consider $H_\tau = 3$ months immediately preceding τ . That is, for every additional month of training or test data that is provided to the model, we use the preceding 3 months to create the historical network and compute the corresponding features on all days in τ . For choosing the experts with an in-degree threshold, we select a threshold of 10 to filter out users having an in-degree of less than 10 in G_{H_τ} from exp_τ . For the centrality features, we set k to be 50; that is, we choose the top 50 users sorted by the corresponding metric in table I. We build different learning models using the GT available from separate *event – types*.

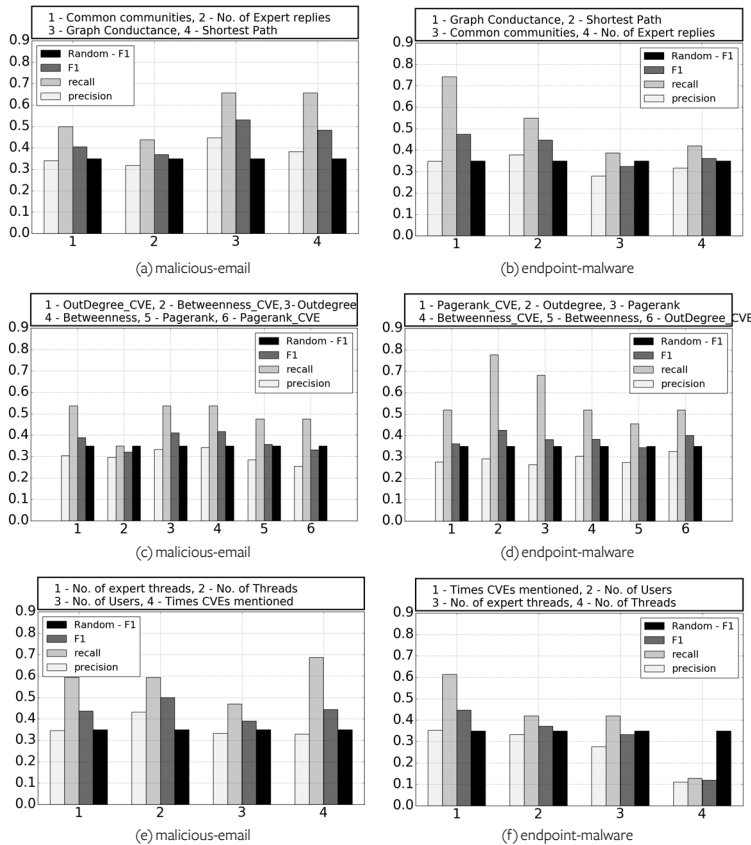


Fig. 4: Classification results for the features considering the logistic regression model: = 7 days, = 8 days.

As mentioned in section III-B, we consider a binary prediction problem in this paper. We assign an attack flag of 1 for at least 1 attack on each day and 0 otherwise. For malicious email, out of the 335 days considered in the dataset, there were reported attacks on 97 days;

this constitutes a positive class ratio of around 29 percent. For endpoint malware, the total number of attack days was 31 out of the 306 days considered in the dataset, which constitutes a positive class ratio of around 26 percent. For evaluating the performance of the models on the dataset, we split the time frame of each event into 70 percent - 30 percent, averaged to the nearest month separately for each event type; that is, we take the first 80 percent of time (in months) as the training dataset and the remaining 20 percent in sequence for the test dataset. We avoid shuffle split as generally being done in cross-validation techniques in order to consider consistency in using sequential information when computing the features. As shown in figure 1, since the period of the attack information provided varies in time for each of the events, we use different time frames for the training model and the test sets. For the event type malicious email, which remains our primary test bed evaluation event, we consider the time period from October 2016 to June 2017 (9 months) in dark web forums for our training data and the period from July 2017 to August 2017 (3 months) as our test dataset. For endpoint malware, we use the time period from April 2016 to September 2016 (6 months) as our training time period and June 2017 to August 2017 (3 months) as our test data for evaluation.

We consider a 1-week time window while keeping $w = 8$ days. From among the set of statistics features that were used for predicting malicious email attacks shown in figure 4(e), we observe the best results using the number of threads as the signal, for which we observe a precision of 0.43, recall of 0.59, and an F1 score of 0.5 against the random F1 of 0.44 for this type of attack. From among the set of expert-centric features in figure 4(a), we obtain the best results from graph conductance with a precision of 0.44, recall of 0.65, and an F1 score of 0.53, which shows an increase in recall over the number of threads measure. Additionally, we observe that the best features in terms of F1 score are graph conductance and shortest paths, whereas the number of threads and vulnerability mentions turns out to be the best among the statistics. For the attacks belonging to the type endpoint malware, we observe similar characteristics for the expert-centric features in figure 4(b), where we obtain a best precision of 0.34, recall of 0.74, and an F1 score of 0.47 against a random F1 of 0.35, followed by the shortest paths measure. However, for the statistical measures, we obtain a precision of 0.35, recall of 0.61, and an F1 score of 0.45 for the vulnerability mentions, followed by the number of threads, which gives us an F1 score of 0.43. Although the common community features do not help much in the overall prediction results, in the following section, we describe a special case that demonstrates the predictive power of the community structure in networks. On the other hand, when we investigate the centrality features with respect to the prediction performance in figure 4(c), we find that just looking at network centralities does not help. The best values we obtain for malicious-email event predictions are from the out-degree and betweenness metrics, both of which give us an F1 score of 0.41. Surprisingly, we find that when the metrics are used for only the users with CVE mentions, the results worse, with the best F1 score for out-degree CVE having an F1 score of 0.38. This calls for a more

complex understanding of path structures between users, rather than just focusing on user significance solely. The challenging nature of the supervised prediction problem is not just due to the issue of class imbalance, but also to the lack of large samples in the dataset which, if they had been present, could have been used for sampling purposes. As an experiment, we also used random forests as the classification model, but we did not observe any significant improvements in the results over the random case.

For the model with the group lasso regularization in equation 2, we set the parameters m , l , and g , and 0.3, 0.3, and 0.1, respectively. We obtained better results for each group of features together for the malicious email event type, with an F1 score of 0.55 for expert-centric, 0.51 for forum/user statistics, and 0.49 for network centrality-based features..

Prediction in High-Activity Weeks

One of the main challenges in predicting external threats without any method for correlating them with external data sources like the dark web or any other database is that it is difficult to validate which kinds of attacks are most correlated with these data sources. To this end, we examine a controlled experiment setup for the malicious email attacks in which we only consider the weeks which exhibited a high frequency of attacks compared to the overall time frame. In our case, we consider weeks that had more than five attacks in the test time frame. These high numbers may be due to multiple attacks in 1 or a few specific days or a few attacks on all days. We run the same supervised prediction method but evaluate them only on these specific weeks.

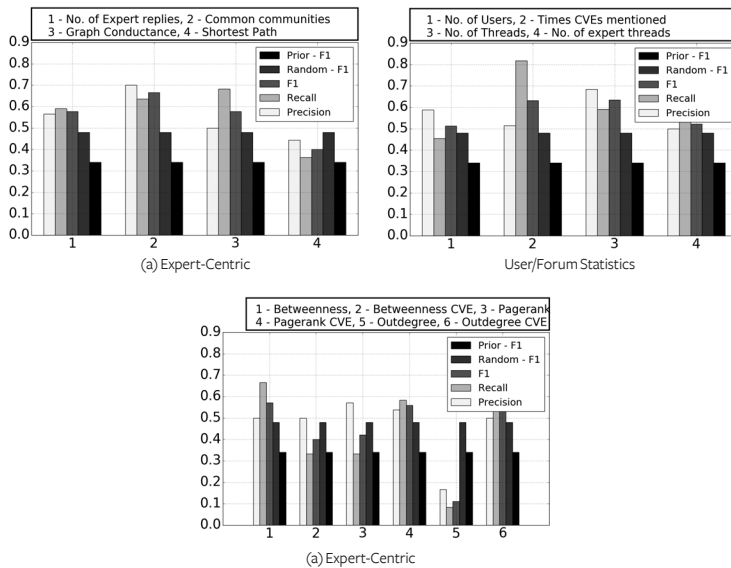


Fig. 5: Classification results for *malicious-email* attacks in high-frequency weeks, $\partial = 7$ days and $n = 8$ days.

From the results shown in figure 5, we find that the best results were shown by the common communities feature, which had a precision of 0.7, a recall of 0.63, and an F1 score of 0.67, compared to the random (no priors) F1 score of 0.48 and a random (with priors) F1 score of 0.34 for the same time parameters. Among the statistics measures, the highest F1 score was 0.63 for the vulnerability mentions feature. From among the set of centrality features, we find that the betweenness measure had the best F1 score (0.58), with a precision of 0.5 and a recall of 0.78. This also suggests the fact that analyzing the path structure between nodes is useful since betweenness relies on the paths passing through a node. Additionally we find that, unlike the results over all of the days, for these specific weeks, the model achieves high precision while maintaining comparable recall, emphasizing the fact that the number of false positive was also reduced during these periods. This correlation between the weeks that exhibit huge attacks and the prediction results imply that network structure analytics can definitely help generate alerts for cyberattacks.

V. RELATED WORK AND CONCLUSION

Using network analysis to understand the topology of dark web forums was studied at breadth in ^[6], where the authors used social network analysis techniques on the reply networks of forums. There have been several attempts to use external social media data sources to predict real-world cyberattacks^{[2],[7],[8]}. The use of machine learning models to predict security threats ^[2] presents many research opportunities, including predicting whether a vulnerability would be exploited based on dark web sources^{[3],[9]}. The availability of large external data sources makes the use of machine learning methods to predict cyberattacks more promising. Previous studies also included the use of time series models to forecast the number of cyber incidents^[16], which increases the demand for the use of such models in cyberattack prediction. The authors in^[17] look at text-mining techniques to understand the content of the posts which provide threat intelligence on various social media platforms. In this study, we argue that the dark web can be a reliable source of information for predicting external enterprise threats. We leverage the network and interaction patterns in the forums to understand the extent to which they can be used as useful indicators. Our study also opens further research possibilities surrounding sentiment analysis on these discussions, which could help track malicious discussions and hence defend against cyber conflict during competition.🛡️

NOTES

1. Samtani, Sagar, Ryan Chinn, and Hsinchun Chen. "Exploring hacker assets in underground forums." IEEE (ISI), 2015.
2. Liu, Yang, et al. "Cloudy with a Chance of Breach: Forecasting Cyber Security Incidents." USENIX Security Symposium. 2015.
3. Nunes, Eric, et al. "Darknet and deepnet mining for proactive cybersecurity threat intelligence." IEEE ISI (2016).
4. Xu, Tingyang, Jiangwen Sun, and Jinbo Bi. "Longitudinal lasso: Jointly learning features and temporal contingency for outcome prediction." ACM, KDD 2015.
5. Meier, Lukas, Sara Van De Geer, and Peter Bühlmann. "The group lasso for logistic regression." Journal of the Royal Statistical Society: Series B (Statistical Methodology) 70.1 (2008): 53-71.
6. Almukaynizi, Mohammed, et al. "Predicting cyber threats through the dynamics of user connectivity in darkweb and deepweb forums." ACM Computational Social Science. (2017).
7. Liu, Yang, et al. "Predicting cyber security incidents using featurebased characterization of network-level malicious activities." 2015 ACM International Workshop Security and Privacy Analytics.
8. Khandpur, Rupinder Paul, et al. "Crowdsourcing cybersecurity: Cyber attack detection using social media." ACM CIKM 2017.
9. Almukaynizi, Mohammed, et al. "Proactive identification of exploits in the wild through vulnerability mentions online." IEEE CyCON, 2017.
10. Thonnard, Olivier, et al. "Are you at risk? Profiling organizations and individuals subject to targeted attacks." International Conference on Financial Cryptography and Data Security. Springer 2015.
11. Bilge, Leyla, and Tudor Dumitras. "Before we knew it: an empirical study of zero-day attacks in the real world." Proceedings of the 2012 ACM conference on Computer and communications security.
12. Sabottke, Carl, Octavian Suci, and Tudor Dumitras. "Vulnerability Disclosure in the Age of Social Media: Exploiting Twitter for Predicting Real-World Exploits." USENIX Security Symposium. 2015.
13. Herley, Cormac, and Dinei Florêncio. "Nobody sells gold for the price of silver: Dishonesty, uncertainty and the underground economy." Economics of information security and privacy. Springer, Boston, MA, 2010. 33-53.
14. Allodi, Luca, Marco Corradin, and Fabio Massacci. "Then and now: On the maturity of the cybercrime markets the lesson that blackhat marketeers learned." IEEE Transactions on Emerging Topics in Computing 4.1 (2016): 35-46.
15. Chen, Hsinchun. "Sentiment and affect analysis of dark web forums: Measuring radicalization on the internet." IEEE ISI 2008.
16. Okutan, Ahmet, et al. "POSTER: Cyber Attack Prediction of Threats from Unconventional Resources (CAPTURE)." Proceedings of the 2017 ACM SIGSAC.
17. Sapienza, Anna, et al. "Early warnings of cyber threats in online discussions." Data Mining Workshops (ICDMW), 2017.