PCAM: A Data-driven Probabilistic Cyber-alert Management Framework

HAIPENG CHEN, Harvard University, USA ANDREW DUNCKLEE, Clark University, USA SUSHIL JAJODIA, George Mason University, USA RUI LIU and SEAN MCNAMARA, Dartmouth College, USA V. S. SUBRAHMANIAN, Northwestern University, USA

We propose PCAM, a Probabilistic Cyber-Alert Management framework, that enables chief information security officers to better manage cyber-alerts. Workers in Cyber Security Operation Centers usually work in 8- or 12-hour shifts. Before a shift, PCAM analyzes data about all past alerts and true alerts during the shift time-frame to schedule a given set of analysts in accordance with workplace constraints so that the expected number of "uncovered" true alerts (i.e., true alerts not shown to an analyst) is minimized. PCAM achieves this by formulating the problem as a bi-level non-linear optimization problem and then shows how to linearize and solve this complex problem. We have tested PCAM extensively. Using statistics derived from 44 days of real-world alert data, we are able to minimize the expected number of true alerts that are not manually examined by a team consisting of junior, senior, and principal analysts. We are also able to identify the optimal mix of junior, senior, and principal analysts needed during both day and night shifts given a budget, outperforming some reasonable baselines. We tested PCAM's proposed schedule (from statistics on 44 days) on a further 6 days of data, using an off-the-shelf false alarm classifier to predict which alerts are real and which ones are false. Moreover, we show experimentally that PCAM is robust to various kinds of errors in the statistics used.

CCS Concepts: • Security and privacy \rightarrow Intrusion/anomaly detection and malware mitigation; • Mathematics of computing \rightarrow Mathematical optimization; • Computing methodologies \rightarrow Artificial intelligence;

Additional Key Words and Phrases: Cyber-Alert Management, Scheduling, Optimization, Prediction

Authors listed in alphabetical order. Work by all authors (except Jajodia) was primarily done while they were at Dartmouth College.

Authors' Addresses: H. Chen, Center for Research on Computation and Society, Department of Computer Science, Harvard University, Boston, MA; email: hpchen@seas.harvard.edu; A. Duncklee, Information Technology & Consulting, Clark University, Wooster, MA; email: aduncklee@clarku.edu; S. Jajodia, Center for Secure Information Systems, George Mason University, Fairfax, VA; email: jajodia@gmu.edu; R. Liu, Department of Computer Science and Institute for Security, Technology, and Society, Dartmouth College, Hanover, NH; email: Rui.Liu.GR@dartmouth.edu; S. McNamara, Information Technology & Consulting, Dartmouth College, Hanover, NH; email: Sean.R.McNamara@dartmouth.edu; V. S. Subrahmanian, Department of Computer Science and Buffet Institute for Global Affairs, Northwestern University, Evanston, IL; email: vss@northwestern.edu.

© 2022 Association for Computing Machinery.

1533-5399/2022/01-ART67 \$15.00 https://doi.org/10.1145/3511101

Parts of this work were funded by ONR grants N00014-18-1-2670, N00014-16-1-2896, and N00014-20-1-2407 and ARO grant W911NF-13-1-0421.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ACM Reference format:

Haipeng Chen, Andrew Duncklee, Sushil Jajodia, Rui Liu, Sean McNamara, and V. S. Subrahmanian. 2022. PCAM: A Data-driven Probabilistic Cyber-alert Management Framework. *ACM Trans. Internet Technol.* 22, 3, Article 67 (January 2022), 24 pages.

https://doi.org/10.1145/3511101

1 INTRODUCTION

Most **Cyber Security Operations Centers (CSOCs)** are flooded by alerts. For instance, while speaking about the Sony breach in 2015, Reference [13] says that "while the tools were able to identify the malicious activity, those alerts were lost in a sea of 40,000 other alerts that same month." Other sources state that "the security operations center (SOC) is drowning in cybersecurity alerts"¹; they go on to state that banks see over 100,000 alerts per day.

We propose **Probabilistic Cyber-Alert Management (PCAM)**, a framework that uses statistics to manage this flood of alerts, while satisfying workplace requirements within a given work shift. Figure 1 shows the architecture of PCAM. (i) As seen on the left side of Figure 1, before a shift starts, PCAM uses historical statistics to compute a schedule of work specifying the analysts' lunch and other breaks. Because each "true alert" must be examined by an analyst, PCAM tries to ensure that this schedule minimizes the expected number of "uncovered" true alerts. (ii) Then, as shown on the right of the figure, during a shift, security products generate alerts. PCAM dynamically uses an off the shelf classifier [6, 10, 14, 20] to predict if an alert is likely to be real ("true" alert) or not ("false" alert). "True" alerts are dynamically assigned to analysts for manual inspection in accordance with the probability (of being a true alert) returned by the classifier. False alarms are discarded.

The organization of the article and the main contributions are as follows. Section 2 describes explains the differences between past work and ours. Section 3 describes the data used by PCAM and provides a high-level overview of how PCAM works. Section 4 describes the technical details behind PCAM. We formalize the problem of maximizing expected true alert coverage and show how it can be expressed via a non-linear bi-level integer programming. Section 5 addresses two major problems that make the above integer program challenging: Nonlinear integer programs are hard to solve and so are bilevel integer programs. We show how to linearize the problem and then additionally show how to modify the bilevel optimization into a single-level optimization. This makes the problem amenable to solution by off the shelf integer program solvers. Section 7 describes our experimental validation of PCAM. We used 30 days of real-world data from Dartmouth College's Information Technology and Consulting (ITC) organization that maintains all of Dartmouth's networks and servers. We then used the real-world alert data for the next 14 days to see how our system would perform if the schedule for the next shift was created at the end of the previous shift. In addition, we ran robustness experiments using simulated data whose statistics varied from those learned from the prior history. In all these cases, PCAM performed extremely well.

2 RELATED WORK

In this section, we introduce related works that study *security personnel scheduling for manual analysis of security alarms.*

The idea that managed security services need to optimize the allocation of analysts was studied in Reference [11]; the paper developed a combination of game theory and probabilistic temporal

¹https://bricata.com/blog/how-many-daily-cybersecurity-alerts/.

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.



Fig. 1. PCAM architecture.

logic to reason about scheduling analysts in the presence of an adversary who takes advantage of the fact that lots of security alerts might be false. The problem was addressed by Dunstatter et al. [4, 5] from a Markov games perspective using deep reinforcement learning. Wang et al. [19] presented a survey of game-theoretic methods that have been applied to improve cyber security. They are focused on a game-theoretic setting with an adversary, which is not the focus of this article.

Ganesan et al. [8, 9] introduced the CSOC workforce optimization problem and used a reinforcement learning method to allocate the incoming alerts to security staff. This is extended in Reference [18] into a setting with distributed CSOCs. Their work does not address what happens before a shift.

Shah et al. [17] proposed an optimization framework to allocate cyber alert detection sensors to alert analysts. This was extended in Reference [16], which focused on fairness issues in the sensoranalyst allocation problem. Okimoto et al. [15] proposed a system to optimize the cyber-security systems involved with multiple criteria, e.g., risk (security), surveillance (privacy), and cost. Altner et al. [1] investigated optimizing staffing and shift scheduling decisions given unknown demand on weekly shifts with an on-call mechanism. In two less-related (as they do not work on analyst scheduling) but still relevant works, Franklin et al. [7] proposed a visualization tool to help with the cyber analysts to understand the analytic process and to structure their analytic thinking. Larriva-Novo et al. [12] proposed a dynamic risk management system with the capability of reacting to those rapid changes in the context of the organization with various types of sensors.

In contrast to all of these efforts that do not rely on a mechanism to separate real from false alerts, PCAM combines scheduling with prediction to (i) comes up with an *a priori* (e.g., just before a shift starts) schedule that minimizes the *expected number of uncovered true alerts* (i.e., true alerts that are missed by analysts) given some resource constraints and some workplace constraints, (ii) come up with a real-time mechanism to handle alerts that come in, once the shift schedule has been created, and (iii) considers variant of PCAM involving robustness to various natural or adversarial scenarios. Moreover, unlike the past work, PCAM has been tested on real-world data and shows high-quality performance. In addition, our formulation of the analyst scheduling problem is closer to the realistic setting with more practical constraints (e.g., capacity of analysts,

constraints on break time, etc.) and therefore introduces a more complex mixed integer nonlinear bi-level optimization problem that has not been addressed in previous works.

3 PCAM DATA

Figure 1 shows the proposed architecture of the PCAM system. PCAM takes as input the data coming from alert-generating products used in enterprises. In our experiments, we used alerts coming from two well-known commercial security products used at Dartmouth College. We used two pieces of information from these products.

- (1) Alert data statistics. We obtained 44 days of real alert data.² These data contained 8,119,858 alerts, leading to a daily average of 180,441 alerts per day. From this entire set of alerts, according to Dartmouth's security analysts, there were just 21,916 true alerts, i.e., about 0.27% of the alerts were real, while about 99.73% of the alerts were false alarms. This constitutes an average of 498 true alerts per day, lost in a sea of 180,441 alerts being generated every day. We obtained statistics about the distribution of these alerts (both all alerts and true alerts) within 10-minute windows during the course of any given day. In other words, we broke the day down into 10-minute windows and obtained statistics (mean and standard deviation of the number of alerts per window, mean and standard deviation of the number of true alerts per window) and used these in Section 4 to generate data-driven schedules for analysts for the next shift so that the shift schedules satisfy various workplace requirements. Figure 2(a) and (b) show the temporal distribution of the number of true alerts during the weekdays and weekends. The results show that the distribution of true alerts is very different during weekdays and weekends. During weekdays, the true alerts largely occur during the 12 noon-to-3 pm window, but during weekends, we see three spikes: during the 3 am-to-4 am window, during the 12 noon-to-3pm window (as was the case during weekdays), as well as during the 6 pm-to-9 pm window. All times are U.S. East Coast times.
- (2) Raw Alerts. In addition, once workers are on a shift, new alerts keep coming in from the security products used by the enterprise. In this case, we developed a "True Alert Predictor" to predict which of the alerts coming in in real time were true and which were false. We were able to do this with high accuracy.

As a consequence, our PCAM architecture allows us to automatically schedule analyst shifts in a manner that is consistent with the arrival of both alerts and true alerts in 10-minute windows prior to the start of a shift; during a shift, the PCAM framework can present to the analysts those alerts that are predicted to be true in descending order of the probability of the truth.

4 PCAM ANALYST SHIFT SCHEDULING AS A NONLINEAR BI-LEVEL OPTIMIZATION

In this section, we describe the problem of scheduling an analyst's shift for **Maximal True Alert Coverage (MTAC)**. Past work [1] suggests that most shifts are either 8-hour shifts or 12-hour shifts. Because our PCAM framework is the same irrespective of the duration of the shift, without loss of generality, we assume the shift is 12 hours long. Twelve-hour shifts usually run from 7 am to 7 pm (day shift) and 7 pm to 7 am (night shift).

We assume a shift is divided up into contiguous time slices t_1, \ldots, t_n where each time slice t_j represents a user-selected interval of time that is sufficiently small for their purposes. In our work, we chose the time slices to be 10 minutes long, and hence a shift of 12 hours corresponds to 72 time slices, i.e., n = 72. We assume that the company employs m analysts a_1, \ldots, a_m . Each analyst is a junior analyst, a senior analyst, or a principal analyst, and their capabilities in terms of processing

²An additional 6 days of data was later used for live testing.

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.



Fig. 2. Temporal distribution of true alerts for weekdays and weekends. The x-axis represents the time and y-axis represents the number of true alerts per hour. The shadowed areas represent the standard deviation.

alerts are given below; the statistics below are reproduced from the study in Reference [1]. We assume that analyst a_i is capable of analyzing c_i alerts per time slice in accordance with the above statistics. In addition, we use the following notation:

- *MT* is a constant integer that denotes the maximum number of time slices that an analyst may work during a shift.
- *CT* is a constant integer that denotes the maximum number of contiguous time slices that an analyst may work before taking a break for one time slice.
- *L* is a constant integer that says that each analyst must have a lunch break consisting of at least *L* contiguous time units. (We recognize that the "lunch break" in the night shift is a break for a meal rather than for lunch, but will abuse notation and call it a lunch break).
- *LS*, *LE* where $LS \le LE$ are constant integers that bound the period when each analyst gets their lunch break. For instance, during the day shift, we may set *LS* to 37 (i.e., 11:00–11:10 am) and *LE* to 51 (i.e., 2:20–2:30 pm), and this imposes the constraint that every analyst gets his lunch break during the time interval [*LS*, *LE*]. In particular, we require that each analyst's lunch break be fully contained within the [*LS*, *LE*] interval.
- We use A_j and TA_j as the expected number of alerts (respectively, true alerts) to arrive during time interval t_j . These numbers can be estimated from historical statistics and computed as expected values.

The problem of scheduling shifts is now the problem of allocating analysts to shifts so that the expected number of "uncovered true alerts" is minimized. We first formalize this as a bi-level mixed integer nonlinear optimization problem as shown in Figure 3.

4.1 Constraints

We now explain the constraints in the mixed integer nonlinear optimization formulation provided in Figure 3.

Our formulation uses binary-valued variables $v_{i,j}$. The idea is that $v_{i,j}$ should be set to 1 if analyst a_i is assigned to work during time slice t_j and 0 otherwise, i.e.,

$$v_{i,j} = \begin{cases} 1 & \text{if analyst } a_i \text{ works during time slice } t_j \\ 0 & \text{otherwise} \end{cases}$$

minimize $\sum_{j=1}^{n} \left(\max\{0, TA_j\} \right)$ subject to	$-\Sigma_{i=1}^m v_{i,j} \cdot c_i\}\Big)$		
$\sum_{j=1}^{n} v_{i,j}$	$\leq MT$	$\forall a_i$	(1)
$\sum_{j=h}^{j=CT+h} v_{i,j}$	$\leq CT - 1$	$\forall a_i, \forall h = 1, \dots, n - CT$	(2)
u _{i,h}	$= \Pi_{k=h}^{k=L+h} (1-v_{i,k})$	$\forall a_i, \forall h = 1, \dots, n - CT$	(3)
$\Sigma_{h=LS}^{LE-L} u_{i,h}$	= 1	$\forall a_i$	(4)
v _{i,j}	$\in \{0,1\}$	$\forall a_i, \forall t_j$	(5)
$ $ $u_{i,h}$	$\in \{0,1\}$	$\forall a_i, \forall h = LS, \dots, LE - L$	(6)



Constraint (1). Constraint (1) of Figure 3 sets one constraint for each analyst saying that that analyst cannot work for more than *MT* time slices during any given shift.

Constraint (2). The goal of the second type of constraint is to ensure that each analyst gets a break periodically. More formally, every time interval of length *CT* is required to have at least one break in it. For this, suppose we set $BR_h = [h, CT + h]$ to be an interval of length *CT* starting at time slice t_h . Hence, we have time intervals such as the following:

$$BR_1 = [1, CT + 1]$$

$$BR_2 = [2, CT + 2]$$

$$\dots$$

$$BR_{n-CT} = [n - CT, n].$$

For each such interval $BR_h = [h, CT + h]$ where $1 \le h \le n - CT$, a break is needed for each analyst. Thus, for each $1 \le i \le m$ and each time slice $1 \le h \le n - CT$, we require that $\sum_{j=h}^{j=CT+h} v_{i,j} \le CT - 1$. By having this constraint for each and every time window BR_h of size CT, we ensure that at least one of the $v_{i,j}$'s is set to 0 that would then be the time slice in which the analyst a_i gets a break.

Constraints (3) and (4). The third constraint requires that every analyst gets a lunch break. For each analyst a_i , suppose $LT_h^i = [h, L+h]$. Intuitively, each LT_h^i interval corresponds to the possible lunch periods for analyst a_i that starts at time t_h and goes until time t_{L+h} . We require this *entire* time to be a break for the analyst, i.e., the value of the variable $v_{i,j}$ should be 0 for each $j \in LT_h^i$ if this interval happens to be the analyst's lunch break. To express this idea, we define an intermediate variable $u_{i,h}$ that is set to 1 if analyst a_i 's lunch break starts at time t_h . Constraint (3) is therefore as follows:

$$u_{i,h} = \prod_{k=h}^{k=L+h} (1-\upsilon_{i,k}).$$

If this is indeed analyst a_i 's lunch break, then $v_{i,h} = 0$ for all $k \le h \le L + h$. This means that $(1 - v_{i,j}) = 1$ for all $k \le h \le L + h$, which in turn means that $\prod_{k=h}^{k=L+h} (1 - v_{i,k})$ would equal 1 *if* the interval $LT_h^i = [h, L + h]$ is indeed the analyst's lunch break.

Because the analyst gets one and only one lunch break, Constraint (4) says that one and only one of the possible $u_{i,h}$'s must be set to 1. Thus, Constraints (3) and (4) jointly enforce the fact that

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.

each analyst gets a contiguous lunch break of length L during his shift. But this comes at a cost, because Constraint (3) is nonlinear.

Constraints (5) and (6). These two constraints ensure that the variables $v_{i,j}$ and $u_{i,h}$ are binary variables.

4.2 Objective Function

The formulation of the MTAC problem *assumes* that we have an estimate of the expected number TA_j of true alerts that occur during any given time period t_j , $1 \le j \le n$ — we can see that we do have this information from our Dartmouth College data as shown in Figure 2(a) and (b).

Figure 2(a) and (b) show that the distribution of true alerts during weekdays is dominated (at least at Dartmouth College) by the 12 pm-to-3 pm window. In contrast, there are multiple peaks during the weekends with attacks between the 3 am-to-4 am windows, as well as between the 12 pm-to-3 pm window and the 6 pm-to-9 pm window. Though these statistics only apply to the true alerts at Dartmouth College, most medium and large enterprises can generate similar distributions.

We note that we can estimate TA_j from historical data to simply be the mean of the set $\{TA_j(d)|d \in TD\}$, i.e., if *j* is the 10-minute time slot between 10 am and 10:10 am, then the mean (respectively, standard deviation) is simply the mean (respectively, standard deviation) of the number of true alerts generated during this time window in the historical data.³

The term $v_{i,j} \cdot c_i$ is the number of alerts that analyst a_i can handle (on average) in one time slice. Note that when $v_{i,j} = 0$, i.e., when the analyst is on a break, he or she can handle no alerts at all. Thus, the summation $\sum_{i=1}^{m} v_{i,j} \cdot c_i$ is the total number of alerts that the *m* analysts can handle in one time slice t_j , and hence $TA_j - \sum_{i=1}^{m} v_{i,j} \cdot c_i$ is an estimate of the total number of alerts that would be left "uncovered" in that one time slice t_j . By uncovered, we mean that this is the number of alerts that are not examined/handled by any of the analysts during this time slice. Because this number could be negative, we take the max of this and 0 in the objective function, i.e., the number of uncovered alerts in time slice t_j would be max $(TA_j - \sum_{i=1}^{m} v_{i,j} \cdot c_i)$. Finally, we sum up the number of uncovered alerts over time to get the total number of uncovered alerts during a given shift, i.e., the total number of uncovered alerts during the entire shift consisting of the time slices t_1, \ldots, t_n is

$$\sum_{j=1}^n \left(\max\left\{ 0, TA_j - \sum_{i=1}^m v_{i,j} \cdot c_i \right\} \right).$$

The goal of our objective function is therefore to minimize the number of uncovered alerts across the shift as a whole.

Discussion: We see immediately that the objective function essentially has two levels, where the outer-level optimization minimizes the total number of uncovered true alerts and the inner-level optimization ensures that the number of uncovered true alerts cannot be "negative."

5 EFFICIENT SOLUTION AS A MIXED INTEGER LINEAR PROGRAM

Because the problem of computing shifts involves optimizing an integer program that is bi-level and nonlinear, it is not directly solvable using existing commercial optimization solvers. In the rest of this section, we show how this problem can be neatly encoded as a single-level mixed integer linear optimization problem.

 $^{^{3}}$ If so desired, then we can compute the mean and standard deviation using just that portion of the training data that looks at the last *w* time windows as opposed to all of them.

5.1 Non-Linear Constraints Transformation

The formulation of our problem of maximizing the expected number of true alerts that are covered is nonlinear, which suggests that it can be very expensive to solve it. In this section, we show that we can represent it in a linear form.

To linearize the nonlinear constraints in Equation (3), we replace each of the nonlinear constraints with two linear constraints:

$$u_{ih} \leq 1 - \sum_{k=h}^{L+h} v_{i,k} / M \quad \forall a_i, \forall h = 1, \dots, n - CT,$$

$$(7)$$

$$u_{ih} \geq 1 - \sum_{k=h}^{L+h} v_{i,k} \qquad \forall a_i, \forall h = 1, \dots, n - CT.$$
(8)

Here M > m is any constant that is larger than the total number of analysts.

The following lemma shows that the constraints are expressed by the following equations.

LEMMA 5.1. In the formalized MTAC problem, the constraints expressed in Equations (5)-(8) have the exact same set of solutions as the Constraints (3), (5), and (6).

Constraints (5), (6), (7), and (8) are equivalent to Equations (3), (5), and (6).

PROOF. Since Equations (5) and (6) require both $v'_{i,j}s$ and u_j 's to be binary variables, we only need to prove the equivalence of Equations (7) and (8) and (3).

We now prove that for each analyst a_i and time slice h = 1, ..., n - CT, the same values of $v_{i,k}, k = h, ..., L + h$ will yield the same values of $u_{i,h}$ in both cases. For ease of reference, we denote the $u_{i,h}$'s in Equation (3) and in Equations (7) and (8) as $u_{i,h}^{I}$ and $u_{i,h}^{II}$, respectively. We split the discussion into the following two cases.

Case 1: $v_{i,k} = 0$, $\forall k = h, ..., L + h$. In this case, it is easy to see that $u_{i,h}^I = \prod_{k=h}^{k=L+h} (1 - v_{i,k}) = 1$. At the same time, we have

$$\begin{split} u_{i,h}^{II} &\leq 1 - \sum_{k=h}^{L+h} v_{i,k} / M = 1, \\ u_{i,h}^{II} &\geq 1 - \sum_{k=h}^{L+h} v_{i,k} = 1, \end{split}$$

which implies that $u_{i,h}^{II}=1.$ Thus, $u_{i,h}^{I}=u_{i,h}^{II}.$

Case 2: $v_{i,k} = 1, \exists k = h, ..., L + h$. Denote the sum term $\sum_{k}^{L+h} v_{i,k}$ as σ . In this case, σ would be an integer value between 1 and m: $1 \le \sigma \le m$. For $u_{i,h}^{I}$, we have $u_{i,h}^{I} = \prod_{k=h}^{k=L+h} (1 - v_{i,k}) = 0$. For $u_{i,h}^{II}$, we have

$$u_{i,h}^{II} \le 1 - \sum_{k=h}^{L+h} v_{i,k}/M = 1 - \sigma/M,$$
$$u_{i,h}^{II} \ge 1 - \sum_{k=h}^{L+h} v_{i,k} = 1 - \sigma.$$

Because $1 \le \sigma \le m \le M$, we have $0 < 1 - \sigma/M < 1$. So $1 - \sigma \le u_{i,h}^{II} < 1$. Combining the constraint that $u_{i,h}$'s are binary variables, we have $u_{i,h}^{II} = 0$. Thus, $u_{i,h}^{I} = u_{i,h}^{II}$ still holds.

This result therefore shows a way of linearizing the constraints while preserving the same set of solutions as the solutions of the initial nonlinear integer program.

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.

minimize subject to	Σ'_j	$\sum_{i=1}^{l} w_j$		
	w_j	$\geq TA_j - \Sigma_{i=1}^m v_{i,j} \cdot c_i$	$\forall t_j$	(9)
	w_j	≥ 0	$\forall t_j$	(10)
	$\sum_{j=1}^{n} v_{i,j}$	$\leq MT$	$\forall a_i$	(11)
	$\sum_{j=h}^{j=CT+h} v_{i,j}$	$\leq CT - 1$	$\forall a_i, \forall h = 1, \dots, n - CT$	(12)
	u_{ih}	$\leq 1 - \sum_{k=h}^{L+h} v_{i,k} / M$	$\forall a_i, \forall h = 1, \dots, n - CT$	(13)
	u_{ih}	$\geq 1 - \Sigma_{k=h}^{L+h} v_{i,k}$	$\forall a_i, \forall h = 1, \dots, n - CT$	(14)
	$\Sigma_{h=LS}^{LE-L} u_{i,h}$	= 1	$\forall a_i$	(15)
	$v_{i,j}$	$\in \{0,1\}$	$\forall a_i, \forall t_j$	(16)
	$u_{i,h}$	$\in \{0, 1\}$	$\forall a_i, \forall h = LS, \dots, LE - L$	(17)

Fig. 4. MILP formulation for Shift Scheduling.

5.2 Bi-Level Objective Function Transformation

As pointed out above, the bi-level objective function is also a source of high computational complexity of the MTAC problem. To handle this problem, we introduce an auxiliary variable w_j , where for each time slice t_j we require that $w_j \ge TA_j - \sum_{i=1}^m v_{i,j} \cdot c_i$ and $w_j \ge 0$. We now show that the optima are preserved as well.

LEMMA 5.2. Minimizing the objective in the optimization problem of Figure 3 is equivalent to the following optimization problem:

minimize
$$\sum_{j=1}^{n} w_j,$$
 (18)

subject to
$$w_j \ge TA_j - \sum_{i=1}^m v_{i,j} \cdot c_i \quad \forall t_j,$$
 (19)

$$w_j \ge 0 \qquad \forall t_j.$$
 (20)

PROOF. Let $\mathbf{v}_j = \langle v_{i,j} \rangle$, i = 1, ..., m denote the decision variable (whether they are working or not) for all the analysts *i* at time slice t_j . For each time slice t_j , we define $f_j(\mathbf{v}_j) := \max\{0, TA_j - \sum_{i=1}^m v_{i,j} \cdot c_{i,j}\}$. The objective can be represented as $f(\mathbf{v}_1, ..., \mathbf{v}_n) = \sum_{j=1}^n f_j(\mathbf{v}_j)$.

Because for each time slice t_j there is $w_j \ge TA_j - \sum_{i=1}^m v_{i,j} \cdot c_i$ and $w_j \ge 0$, we thus have $w_j \ge f(\mathbf{v}_j)$. Let $w = \sum_{j=1}^n w_j$; we then have $w \ge f(\mathbf{v}_1, \ldots, \mathbf{v}_n)$. Therefore, we can define the epigraph of function $f(\mathbf{v}_1, \ldots, \mathbf{v}_n)$ as

$$epif = \{(\mathbf{v}, w) | \mathbf{v} \in domf, f(\mathbf{v}) \le w\},\$$

where $\mathbf{v} = \langle \mathbf{v}_j \rangle$, j = 1, ..., n is the decision variable for all the analysts throughout all the time slices. It is easy to see that the optimization described in Equations (18)–(20) is the epigraph form of minimizing the objective in Figure 3. According to Reference [2], the epigraph form of the original optimization problem is equivalent.

Combining Lemmas 5.1–5.2, we immediately have the following theorem.

THEOREM 5.3. The bi-level non-linear formulation of the MTAC problem described in Figure 3 is equivalent to the single-level **Mixed Integer Linear Programming (MILP)** in Figure 4.

Important Note. The above theorem is very important, because it reduces the problem of shift scheduling to maximize expected true alert coverage (or alternatively minimize uncovered true alerts) to a single-level MILP problem instead of a bi-level mixed integer nonlinear programming problem. Though MILPs are in general NP-hard to solve, they have been well studied, and there are existing commercial optimization solvers (e.g., Gurobi and GLPK) that can be used to efficiently solve these types of optimization problems. For instance, in our experiments with a team of 6 junior, 8 senior, and 8 principal analysts, the problem can be solved within 0.4943 seconds (averaged over 50 runs, using Gurobi).

Note on Identifying the Optimal Mix of Personnel. Given a (biweekly) budget *B* for the day/night shift, a CISO (or similar senior leader) of an organization needs to identify the optimal mix of people to hire. This can be easily done. Suppose Sal_j , Sal_s , and Sal_p denote the biweekly budget for salaries of a junior, senior, and principal analyst, respectively. Then the organization can hire up to n_j , n_s , n_p junior, senior, and principal analysts, respectively, where $n_j = \lfloor B/Sal_j \rfloor$, $n_s = \lfloor B/Sal_s \rfloor$, and $n_p = \lfloor B/Sal_p \rfloor$. We can then identify the optimal mix of junior, senior, and principal analysts to be budget to be budget via a simple procedure. A staffing triple is a triple of the form (K_j, K_s, K_p) where $1 \le K_j \le n_j$, $1 \le K_s \le n_s$, and $1 \le K_p \le n_p$. Given a budget *B*, let ST(B) denote the set of all staffing triples. For each staffing triple $st_{j,s,p} = (K_j, K_s, K_p) \in ST(B)$, let MTA_{K_j,K_s,K_p} be the number of missed true alerts computed by running the Mixed Integer Linear Program shown in Figure 4. The optimal staffing is then

$$(J^*, S^*, P^*) = \underset{(K_j, K_s, K_P) \in ST(B)}{\arg \min} MTA_{K_j, K_s, K_p}$$

6 IMPROVING ROBUSTNESS OF SCHEDULING

In the above, we use the average number of true alerts to derive a schedule for the analysts. However, the actual number of true alerts might be different from the averages due to (1) fluctuations/randomness or (2) existence of an adversary. To handle these two cases and therefore improve the robustness of PCAM scheduling, we improve the above scheduling algorithm by taking these two factors into account. Our design is inspired by the idea of "adversarial training" in adversarial machine learning methods, where the idea is to add additional "adversarial samples" in the training set so that the learned models are robust to adversarial attacks.

More specifically, the "adversarial sample" is a given true alert distribution $TA_1^{(l)}, \ldots, TA_{MT}^{(l)}$ denoted by a superscript *l*. The adversarial true alert distribution can be different from the average numbers of the distribution mwhich the probabilities used by PCAM are derived. For each such distribution $TA_1^{(l)}, \ldots, TA_{MT}^{(l)}$, we add it into the constraints in Equation (9),

$$w_j \ge TA_j^{(l)} - \sum_{i=1}^m v_{i,j} \cdot c_i \quad \forall t_j.$$

$$(21)$$

With sufficient sampling of such adversarial distributions, we can obtain a super-set of constraints that take into account changes of true alert distributions. In particular, we consider the following three types of changes.

PCAM-Fluct: In the first case, we consider fluctuations of true alert distributions for each time slice. This is done by sampling from a probability distribution (a Gaussian) with the mean \overline{TA}_j and standard deviation std_j obtained from the historical data:

$$TA_i^{fluct} \sim \mathcal{N}(\overline{TA}_j, \operatorname{std}_j), \quad \forall t_j.$$
 (22)

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.

PCAM: A Data-driven Probabilistic Cyber-alert Management Framework 67:11

PCAM-Shift: In the second case, we consider a "shift" of the true alert distributions by τ time slices:

$$TA_{j}^{shift} = \begin{cases} \overline{TA}_{j+\tau}, & \text{if } j+\tau \leq MT\\ \overline{TA}_{j+\tau-MT}, & \text{if } j+\tau > MT \end{cases}.$$
(23)

Here τ is an integer value that can be positive or negative.

PCAM-Mix: In the last type, we consider a mix of the above two cases, i.e., we first do a shift of the distribution and then perturb it using a Gaussian distribution:

$$TA_j^{mix} = \begin{cases} TA_{j+\tau}^{shift}, & \text{if } j+\tau \le MT\\ TA_{j+\tau-MT}^{shift}, & \text{if } j+\tau > MT \end{cases}.$$
(24)

Note that reversing the fluctuation and shift operations would result in equivalent true alert distributions.

7 EXPERIMENTAL EVALUATION

We now describe the results of our experiments. We only show results of the PCAM Analyst Shift Scheduler in the main text—as predicting whether alerts are real or fake is not a contribution of this article, we show those results in the appendix.

7.1 PCAM Analyst Shift Scheduler

We split 24-hour periods into a day shift (7 am to 7 pm) and a night shift (7 pm to 7 am). We consider time slices of 10 minutes; hence, our shifts have 72 time slices in them.

The average number of historical true alerts during any given day (respectively, night) time slice is used to estimate the number of true alerts during the same time slice in the next day (respectively, night) shift.

7.1.1 Settings. As mentioned earlier, our goal is to come up with a shift schedule that satisfies various workplace constraints while minimizing the number of uncovered true alerts. As the number of uncovered true alerts obviously depends upon the number of analysts at different levels (junior, senior, and principal analysts) that in turn depends on the organization's budget, we vary these numbers in the range $\{2, 3, \ldots, 9\}$. Note that these are used to evaluate how our method works under different settings with different numbers of analysts. Our main formulation is still the bi-level optimization introduced in Figure 3 with fixed number of analysts as part of the optimization.

For each combination tuple of junior, senior, and principal analysts, we use existing optimization solvers such as Gurobi and GLPK to solve the MILP in Figure 4 to optimize (minimize) the number of uncovered true alerts, given their pay and capabilities as described in Table 1. It is important to note that PCAM can of course also be used in settings where different true alert rates hold and where the pay rates and capabilities of analysts are different by simply plugging the new values in instead of those we use in our experiments.

Baseline: We introduce a *baseline* that maximizes the amount of time worked by each analyst, subject to the workplace constraints. Specifically, we use the result of the optimization problem posed in Figure 5 as our baseline with which to compare.

7.1.2 *PCAM Shift Scheduling Results.* Figure 6 compares the PCAM-optimized schedule and the *baseline* for day shifts and night shifts, respectively. The *x*-axis shows the budget. Given a budget, PCAM finds the optimal schedule (i.e., mix of junior, senior and principal analysts) to hire and

	APR	Salary	Glassdoor (NY)	Glassdoor (CA)	Glassdoor (US)
Junior	5	\$3000	\$2500	\$3200	\$2500
Senior	7.5	\$4000	\$2750	\$3600	\$2700
Principal	10	\$6000	\$3400	\$4500	\$3375

Table 1. Biweekly Pay for Analysts

APR denotes Alert Processing Rate/hr.

maximize	$\sum_{i=1}^{n} \sum_{i=1}^{m} v_{i,j}$	
subject to Equation (14) to (20).	

Fig. 5. Baseline shift scheduling.

compares the number of uncovered true alerts during the day and night, respectively. For the day shift, the optimal number found by PCAM to get the number of uncovered true alerts to 0 consists of 6 junior, 8 senior, and 8 principal analysts at a cost of \$98K for a 2-week period. In contrast, for the night shift, PCAM finds that the best combination consists of 3 junior, 7 senior, and 3 principal analysts at a cost of \$55K for a biweekly period.⁴ The reason fewer analysts are needed during the night (at least for the Dartmouth College setting) is intuitively, because there are fewer true alerts during the night on weekdays.

The figures show that the PCAM-optimized schedule beats the baseline by a hefty margin. It is better in two respects:

- (1) PCAM delivers a superior uncovered true alert rate compared to the baseline. At 98K biweekly, PCAM has 0 uncovered true alerts, while the baseline has about 45.
- (2) PCAM saves organizations money—to achieve a day time rate of 0 uncovered true alerts, the baseline would need \$147K for a biweekly period instead of the \$98K needed by PCAM. To achieve a night time rate of 0 uncovered true alerts, the baseline would need \$83K for a biweekly period as compared to the \$55K required by PCAM.

Financial Implications. In short, the use of PCAM would save \$49K per biweekly period for daytime salaries and \$28K per biweekly period for nighttime salaries. As a consequence, the total saving in a year if PCAM was used instead of the baseline would be \$77K ×24 biweekly periods, which is \$1.848M per year. Please note that this estimated saving depends upon many factors— the salaries noted in Reference [1] and the true alert distributions shown in Figure 2(b) and (a). These numbers might vary for other organizations depending upon these factors. Moreover, these numbers assume that the organization in question would like to reduce the expected number of uncovered true alerts to be near zero. If they are willing to take a risk, then the savings might drop. Nonetheless, our methods can be used by any organization as long as they can provide these numbers to PCAM as inputs.

Detailed comparison. Figures 7 and 8 show the detailed number of uncovered true alerts when varying exactly one of the three n_{junior} , n_{senior} , $n_{principal}$, while keeping the other two fixed using both the PCAM-generated schedule and the *Baseline* schedule for day shifts.⁵ As an example,

⁴Please note that the salary rates used were from Reference [1] and from glassdoor.com as opposed to Dartmouth College salaries, which were kept confidential from us. Moreover, benefits are not included in these costs but would be around 30% of salary. Finally, please note that the number of junior, senior, and principal analysts computed above were based on the alert data for Dartmouth College's network. Though these numbers can be different for other organizations, the process followed in this article can be directly applied to the alert statistics from those organizations.

⁵Results for night shifts are similar and we put them in the appendix for space concerns.



(a) For daytime shift, PCAM needs 6 junior, 8 senior, and 8 principal analysts and a budget of 98.0k total pay biweekly is needed.

Night



(b) For nighttime shift, PCAM's solution to achieve 0 uncovered true alerts needs 3 junior, 7 senior, and 3 principal analysts — with a 55.0K total pay biweekly.

0.114

0.009

Fig. 6. Number of uncovered true alerts of optimized schedule and baseline schedule for day and night shifts. The *x*-axes show the biweekly. The *y*-axes show the number of uncovered true alerts. The vertical orange line marks when the uncovered true alerts for the PCAM scheduler is 0.

	crea nat	es under	various	/ losenee	nuces	
Humanistic absence	0	0.2	0.4	0.6	0.8	1.0
Dav	1 000	0 763	0 529	0 310	0 1 3 0	0.00

0.751

0.516

0.295

1.000

Table 2. Uncovered Rates under Various Absence Rates

the first subfigure in Figure 7 shows that when $n_{junior} = 2$, by varying n_{senior} and $n_{principal}$ the uncovered alerts could be reduced to 10.2 using the PCAM-optimized day shift. However, in Figure 8, the subfigure shows that when $n_{junior} = 2$, by varying n_{senior} and $n_{principal}$ the uncovered alerts could be reduced to 58.3 using baseline day shift. The figures show that the PCAM-optimized schedule is substantially better than the *baseline* for every combination of types of analysts hired.

Humanistic practices. From a humanistic perspective, work shifts may be affected by various factors such as sickness, family issues, and so on. We use the **uncovered rate (UR)** of true alerts to evaluate the performance, which is defined as the number of uncovered true alerts divided by the number of total true alerts. We define the humanistic absence rate to to be the absent work capacity divided by the total work capacity. Table 2 shows the uncovered rates for the alerts under various humanistic absence rates. The uncovered rate is obtained by the simulations based on historical data. The table shows that PCAM is able to cover more than 75% of the alerts when the humanistic absence is 0.2 for both day and night simulations. A 20% absentee rate is quite high. The table therefore shows that even when the absentee rate of analysts is high, PCAM still achieves good performance and is robust to the humanistic absence of analysts.

7.2 PCAM Shift Schedule Performance with Adversary

We also tested the performance of PCAM in the presence of an adversary. For each type of adversarial distribution, we also report the *Euclidean Distance*⁶ from the original distribution. We use

⁶We also consider other distance metrics such as Bray–Curtis Distance, Cosine Distance and Wave–Hedges Distance. We refer to Reference [3] for the definitions of the distance metrics. The results using these other distributions are similar and are shown in the appendix.

67:14



Fig. 7. Results for *PCAM* Optimized day shift. Each plot shows the number of uncovered alerts when varying one number of the combination tuple (junior, senior, or principal), while the other two are fixed (the fixed values are shown in the titles of the subfigures). The first, second, and third rows, respectively, fix the number of juniors, senior, or principal to hire. The columns indicate the number of the type of analysts that are fixed. The *x*- and *y*-axes represent the number of the two types of analysts that are being evaluated. The *z*-axis represents the uncovered alerts.

the UR of true alerts to evaluate the performance, which is defined as the number of uncovered true alerts divided by the number of total true alerts.

Fluctuations of true alert distributions: In this experiment, we consider an adversarial setting where the true alert distributions are made to fluctuate around the average values of the historical data. To do this, we randomly sample 100 true alert distributions using the Normal distribution in Equation (22) with a confidence interval (0.95). Table 3 shows the performance of PCAM-Fluct and PCAM in this setting. We immediately see that by taking the fluctuations into account in the formulated optimization problem, PCAM-Fluct performs much better than the vanilla PCAM for both day and night schedules. Figure 9(a) and (d) show the performance of PCAM-Fluct with respect to the Euclidean distance between the original and adversarial true alert distributions. The figures show that the uncovered rate UR increases as the distance between the distributions increases. This is intuitive, as attacks that are more different from the original distribution will have higher success rate.

Shifts of true alert distributions: In this experiment, the adversarial samples are generated by shifting the original true alert distribution by 0.5, 1, 1.5, ..., 6 (i.e., with 12 types of shift periods). Table 4 shows the performance of PCAM-Shift and PCAM in this scenario. As in the preceding experiment, we see that the use of adversarial training enables PCAM-Shift to outperform vanilla PCAM. Figure 9(b) and (e) show the performance of PCAM-Shift as the lengths of the time shift

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.



Fig. 8. Results for *baseline* day shift. Each plot shows the number of uncovered alerts when varying one number of the combination tuple (junior, senior, or principal), while the other two are fixed (the fixed values are in titles of the subfigures). The meanings of the rows, columns, and x-, y-, and z-axes are the same as in Figure 7.

	Total true	UR_{PCAM}	UR _{PCAM-Fluct}	Distance
Day	502.270	0.717	0.340	132.437
Night	498.220	0.772	0.416	154.811

Table 3. Performance of PCAM-Fluct and PCAM with Fluctuations in True Alert Distribution

Results are averaged from 100 samples for both day and night schedules.

increases. In this setting, we can see that there is no clear pattern of UR with respect to the length of the temporal shifts. Our conjecture is that the true alert distribution with larger hours of shift to the original distribution is not necessarily further away from the original distribution.

Mixed adversarial true alert distribution: In this experiment, we consider both fluctuations and time shifts in the true alert distribution. Table 5 shows the performance of PCAM-Mix and PCAM in this scenario. We can see that PCAM-Mix consistently outperforms PCAM by a large margin. Figure 9(c) and (f) show the performance of PCAM-Mix with respect to the distance to the original true alert distribution. As in the case of PCAM-Fluct, the figures show that the Uncovered Rate increases as the distance increases.

7.3 Live test of PCAM

The acid test of the PCAM architecture would be to see how the entire PCAM architecture holds up. If analysts were scheduled in accordance with the schedule computed by PCAM and then the

H. Chen et al.



Fig. 9. Performance of PCAM-Fluct, PCAM-Shift, and PCAM-Mix with respect to distance of the original and adversarial true alert distributions, for day (upper row) and night (lower row) schedules. The *x*-axis represents the *Euclidean Distance* or number of shifts. The *y*-axis represents the UR by each schedule. Results of other distance metrics are shown in the appendix.

Table 4. Performance of PCAM-Shift and PCAM with Shift in True Alert Distribution

	Total true	UR _{PCAM}	UR _{PCAM-Shift}	Distance
Day	279.000	0.434	0.021	58.348
Night	213.000	0.425	0.033	42.795

Results are averaged for 12 samples for both day and night schedules.

Table 5.	Performance	of PCA	M-Mix and	PCAM wit	h Both
Flue	ctuations and	Shifts ir	n True Aler	t Distributio	on

	Total true	UR_{PCAM}	$UR_{PCAM-Mix}$	Distance
Day	502.270	0.745	0.340	143.456
Night	498.220	0.801	0.409	166.163

Results are averaged for 100 samples for both day and night schedules.

real-world alerts occurred as they did occur in the live test, then how many true alerts would be uncovered.⁷

We also did an end-to-end test using the previously determined numbers of analysts for the day and night shifts respectively on the live test (6 days). In this live test, we created the analyst schedules before the day and night shifts and then played back the actual arrival of alerts during the 6 days of the testing period. *Not a single true alert went uncovered by the schedules generated by the PCAM model.*

⁷The schedule assumed 6 junior, 8 senior, and 8 principal analysts for day time schedules and 3 junior, 7 senior, and 3 principal analysts for night time schedules.

ACM Transactions on Internet Technology, Vol. 22, No. 3, Article 67. Publication date: January 2022.

PCAM: A Data-driven Probabilistic Cyber-alert Management Framework

8 CONCLUSION

Most cyber-security operations centers experience a huge flood of security alerts. This deluge is typically managed by a small team that is usually overworked. Because of the sheer volume of alerts, true alerts can be missed from a sea of false alarms.

In this article, we propose a data driven bi-level optimization algorithm to solve the following problems. First, given a distribution of true alerts during the day, how best should a given set of analysts be scheduled to minimize the expected number of true alerts that are not handled by an analyst. Second, given a desire to reduce the expected number of true alerts that are uncovered to be below a desired upper bound, how many analysts of different types do we need? Third, how can a predictor that predicts whether a given alert is real or merely a false alarm be incorporated into the system?

To solve these problems, we propose the PCAM system. PCAM has a "before the shift" mode in which distributions of true alerts over time are used to create a schedule, and a "during the shift" mode in which alerts predicted to be true by PCAM are handled by analysts who are working during that time period. We show that PCAM's end-to-end performance is excellent. Using a 44-day training period with real alert data to train the PCAM model and to capture the desired distributions of true alerts during the day, we are able to show that the end to end PCAM system created schedules that did not miss a single true alert during a 6-day test window. Moreover, our predictive model also generated excellent results during live testing over the 6-day window, outperforming the performance in training. Last, we also show that PCAM does well under various kinds of adversarial (or natural) fluctuations/shifts between the statistics used during training and the statistics seen when the system is subsequently used.

APPENDICES

A PCAM ALERT CLASSIFIER

We also briefly introduce our alert classifier that distinguishes true alerts from false ones. A number of threat detection systems [6, 10, 14, 20] have sought to reduce the false alarm rate. As we can use any such framework within PCAM, we do not discuss this further here. Note that we are not claiming technical novelty in this component but describing it nonetheless for completeness of the PCAM architecture. We used the output logs of two security products used at Dartmouth College's ITC offices, who run Dartmouth's network to define three broad categories of features related to alerts. These include the following:

- *Systems-Related Features:* Systems related features included information such as the time of receipt time (of a suspect packet), the identity of the logging server, and more.
- *Traffic-Related Features:* Traffic related features contain proxies for IP addresses, communication ports, the transmission protocol used, and more.
- *Threat-Related Features:* Threat-related features include the category of the threat and the severity of the threat according to that security product.

We preprocessed data by categorizing sparse attributes, discarding single-value attributes, duplicate attributes, and attributes with too many missing values. Table 6 describes the different types of variables we finally ended up with.

We ran 10 off the shelf classifiers to learn classifiers that separate true alerts from false ones. Seven of the 10 classifiers are "traditional" classifiers: Decision Trees, Logistic Regression, Bernoulli Naive Bayes, Gaussian Naive Bayes, Multinomial Naive Bayes, Random Forest, and **Support Vector Machines (SVM)**. The other three are neural classifiers including **Multi-Layer Perceptrons (MLP)**, **Google's Deep and Wide system (DeepWide)** and **Convolutional Neural Nets**

Serial Number	Serial number of the firewall that generated the log.
Log Forwarding Profile	The profile used for log filtering and tagging.
Source Zone	Indicator of the direction of the connection, which is either outgoing of incom-
	ing.
Ingress Interface	The network interface device that the session was sourced form.
Egress Interface	The network interface device interface that the session was destined to.
Repeat Count	Number of sessions with same Source IP, Destination IP, Application, and Sub-
	type seen within 5 seconds.
Flags	32-bit field that provides details on session, including PCAP (packet capture),
	IPv6, SSL, and so on.
Protocol	IP protocol associated with the session.
Direction	Indicates the direction of the attack, client-to-server or server-to-client.
Source Location	Source country or Internal region for private addresses.
Destination Location	Destination country or Internal region for private addresses.
Source Port	Source port utilized by the session.
NAT Source Port	The source port used by the connection after NAT (network address translation).
Destination Port	Destination port utilized by the session.
NAT Destination Port	The destination port used by the connection after NAT (network address trans-
	lation).
Subtype	Subtype of the threat log. Vulnerability means this threat is vulnerability exploit
	detected via a Vulnerability Protection profile.
Action	Action taken for the session.
Category	The category of the URL.
Severity	Severity associated with the threat. Reported by threat detectors.
Rule Name	Name of the rule that the session matched.

Table 6. Details of Variables

(CNNs). We do not make any claim of novelty about these algorithms—as they give very high F1-scores and AUCs, there was no point developing a new classifier.

B TRUE ALERT PREDICTION RESULTS

The shift schedules experiments above are used to schedule each analyst's scheduled breaks and lunch slots—but they are finalized just *before* the analyst starts his shift.

The True Alert Predictor component of PCAM looks at alerts as they come in during a shift and tries to classify each alert as true or false when the alert comes in. We ran two experiments to assess the accuracy of the True Alert Predictor. In the first experiment, we did a traditional ML analysis of 44 days of data collected by Dartmouth's ITC staff. In the second experiment, we did a 6-day live test of the models learned on the 44 days of data to see how well they performed. As this component only uses traditional classifiers and is not particularly novel, we do not claim this as a contribution of this article—just a necessary part of the PCAM architecture.

Historical Data Analysis. We tested each of the seven classifiers on the features described in Section A. Because our alerts reflect temporal data, cross validation is not an appropriate test.⁸ To avoid the pitfalls of using cross validation for temporal data, we use *rolling window prediction*.

⁸This is because *K*-fold cross validation would *randomly* split the data into *K* chunks. It would then do *K* iterations: In each iteration, one of the *K* chunks would be the test set, while the remaining (K-1) chunks would be used for training the model. The performance of a classifier would then be obtained by aggregating the performance over the *K* folds. However, this is inappropriate for temporally sensitive data, because the training folds can (with $\frac{K-1}{K}$ probability) contain data points from the future and those might be used to predict outcomes for data in the test fold that might be from the past.

	ROC-AUC	F1 score	Precision	Recall
Decision tree	0.9936	0.8634	0.8866	0.8595
Logistic	0.9992	0.8449	0.8918	0.8124
NB Bernoulli	0.9970	0.2621	0.1572	0.9064
NB Gaussian	0.9905	0.2507	0.1471	0.9941
NB multinomial	0.9967	0.2713	0.1639	0.9059
Random forest	0.9957	0.8806	0.9073	0.8600
SVM	0.9967	0.8538	0.9258	0.8064
DNN-MLP	0.9998	0.8740	0.9350	0.8300
DNN-DeepWide	0.9856	0.8634	0.8990	0.8466
DNN-CNN	0.9810	0.8634	0.8833	0.8300

Table 7. Classifier Performance (True Alert vs. False Alert prediction) on 44 Days of Data Using Rolling Window Prediction

Table 8. Live Test: Performance of Random Forest Classifier on True Alert vs. False Alert Prediction on 44 days of live data. Random Forest has the best F1 among all the classifiers

	ROC-AUC	F1 score	Precision	Recall
Day 1	0.9931	0.9677	0.9633	0.9722
Day 2	0.9786	0.9649	0.9910	0.9402
Day 3	0.9998	0.9672	0.9987	0.9365
Day 4	0.9913	0.9652	0.9652	0.9652
Day 5	0.9758	0.9374	0.9450	0.9299
Day 6	0.9867	0.9364	0.9717	0.9035

In rolling window prediction, we train our models for the first *T* days and then predict for day T + 1—then we learn the model from the first T + 1 days of data and predict for day T + 2, and so forth. We now report the results for rolling window prediction—for the sake of completeness, Table 7 contains the results of 10-fold cross validation as well. In our rolling window experiments, we train on each of the intervals $[1, 30], [1, 31], \ldots, [1, 43]$ and then predict for days $31, 32, \ldots, 44$ respectively, i.e., we train on the first *T* days and then make predictions for day T + 1 with $T = 30, \ldots, 43$.

Table 7 shows that PCAM's performance on rolling window prediction is excellent, with Random Forest yielding the best F1-Score of 88.06% with a precision of 90.73% and a recall of 86%. In terms of AUC, the results were almost a dead heat between Random Forest and SVM, with Google's DeepWide coming out very marginally ahead. We therefore feel that Random Forest was the best of the 10 classifiers tested.⁹

Live Testing Experiments. We also conducted a live test. In this test, we used the best model (Random Forest) obtained via training on the 44-day window—but then tested it on 6 days of live data to assess the predictive efficacy of our models. Table 8 shows the result of rolling window prediction in this case.

 $^{^{9}}$ In contrast, the results of K-fold cross validation shown in Table 9 are slightly inflated (due to the methodological flaw with using cross validation on temporal data) with Random Forest yielding the best F1-Score of 90.85% with a precision of 95.19% and a recall of 86.88%.

	ROC-AUC	F1 score	Precision	Recall
Decision tree	0.9895	0.9065	0.9488	0.8678
Logistic	0.9992	0.8898	0.9455	0.8404
NB Bernoulli	0.9967	0.2537	0.1476	0.9018
NB Gaussian	0.9906	0.2510	0.1437	0.9915
NB multinomial	0.9968	0.2643	0.1548	0.9013
Random forest	0.9923	0.9085	0.9519	0.8688
SVM	0.9974	0.8887	0.9488	0.8357
DNN-MLP	0.9998	0.9044	0.9614	0.8539
DNN-DeepWide	0.9825	0.8616	0.9576	0.7830
DNN-CNN	0.9840	0.8765	0.9256	0.8378

Table 9. Experimental Results Showing Classifier Performance in Predicting Whether an Alert Is Real or False Using Cross Validation

We see that the predictive accuracy of Random Forest actually *increased*: the F1-scores for each of the days went up to 93.64—96.77% with precision also increasing to lie between 94.5% and 99.87% and recall to lie between 90.35% and 97.22%. The fact that the trained models *improved* during live testing suggests that over-fitting did not occur during the training phase.

C PCAM DETAILED SCHEDULING RESULTS FOR NIGHT SHIFTS

Figures 10 and 11 show the detailed number of uncovered true alerts when varying exactly one of the three n_{junior} , n_{senior} , $n_{principal}$, while keeping the other two fixed using both the PCAM-generated schedule and the *Baseline* schedule for night shifts. Results show that the PCAM-optimized schedule is substantially better than the *baseline* for every combination of types of analysts hired during night shifts.



Fig. 10. Results for *PCAM* Optimized night shift. Each plot shows the number of uncovered alerts when varying one number of the combination tuple (junior, senior, or principal), while the other two are fixed. The meanings of the rows, columns, and x-, y-, and z-axes are the same as in Figure 7.



Fig. 11. Results for *baseline* night shift. Each plot shows the number of uncovered alerts when varying one number of the combination tuple (junior, senior, or principal), while the other two are fixed. The meanings of the rows, columns, and x-, y-, and z-axes are the same as in Figure 7.



(a) PCAM-Fluct for day with *Euclidean Distance*.



(e) PCAM-Fluct for night with *Euclidean Distance.*



(i) PCAM-Mix for day with *Euclidean Distance*.



(m) PCAM-Mix for night with *Euclidean Distance.*



(b) PCAM-Fluct for day with *Bray-Curtis Distance*.



(f) PCAM-Fluct for night with *Bray-Curtis Distance*.



(j) PCAM-Mix for day with *Bray-Curtis Distance*.



(n) PCAM-Mix for night with *Bray-Curtis Distance*.



(c) PCAM-Fluct for day with *Cosine Distance*.



(g) PCAM-Fluct for nights with *Cosine Distance.*



(k) PCAM-Mix for day with *Cosine Distance*.



(o) PCAM-Mix for night

with Cosine Distance.

(d) PCAM-Fluct for day with *Wave-Hedges Distance*.



(h) PCAM-Fluct for night with *Wave-Hedges Distance*.



(I) PCAM-Mix for day with *Wave-Hedges Distance*.



(p) PCAM-Mix for night with *Wave-Hedges Distance*.

Fig. 12. Performance for day and night schedules of PCAM-Fluct and PCAM-Mix. The x-axis represents the distance metrics, including *Euclidean Distance*, *Bray–Curtis Distance*, *Cosine Distance*, and *Wave–Hedges Distance*. The y-axis represents the Uncovered Rate by each schedule.

REFERENCES

- Douglas S. Altner, Anthony C. Rojas, and Leslie D. Servi. 2018. A two-stage stochastic program for multi-shift, multianalyst, workforce optimization with multiple on-call options. J. Schedul. 21, 5 (2018), 517–531.
- [2] Stephen Boyd and Lieven Vandenberghe. 2004. Convex Optimization. Cambridge University Press.
- [3] Sung-Hyuk Cha. 2007. Comprehensive survey on distance/similarity measures between probability density functions. *City* 1, 2 (2007), 1.
- [4] Noah Dunstatter, Mina Guirguis, and Alireza Tahsini. 2018. Allocating security analysts to cyber alerts using markov games. In Proceedings of the National Cyber Summit (NCS'18). IEEE, 16–23.
- [5] Noah Dunstatter, Alireza Tahsini, Mina Guirguis, and Jelena Tešić. 2019. Solving cyber alert allocation markov games with deep reinforcement learning. In Proceedings of the International Conference on Decision and Game Theory for Security. Springer, 164–183.

- [6] Meisam Eslahi, Habibah Hashim, and Nooritawati Tahir. 2013. An efficient false alarm reduction approach in HTTPbased botnet detection. In Proceedings of the IEEE Symposium on Computers & Informatics (ISCI'13). 201–205.
- [7] Lyndsey Franklin, Meg Pirrung, Leslie Blaha, Michelle Dowling, and Mi Feng. 2017. Toward a visualization-supported workflow for cyber alert management using threat models and human-centered design. In Proceedings of the IEEE Symposium on Visualization for Cyber Security (VizSec'17). IEEE, 1–8.
- [8] Rajesh Ganesan, Sushil Jajodia, and Hasan Cam. 2017. Optimal scheduling of cybersecurity analysts for minimizing risk. ACM Trans. Intell. Syst. Technol. 8, 4 (2017), 52.
- [9] Rajesh Ganesan, Ankit Shah, Sushil Jajodia, and Hasan Cam. 2019. Optimizing alert data management processes at a cyber security operations center. In Adversarial and Uncertain Reasoning for Adaptive Cyber Defense. Springer, 206–231.
- [10] Neminath Hubballi and Vinoth Suryanarayanan. 2014. False alarm minimization techniques in signature-based intrusion detection systems: A survey. Comput. Commun. 49 (2014), 1–17.
- [11] Sushil Jajodia, Noseong Park, Edoardo Serra, and V. S. Subrahmanian. 2016. Using temporal probabilistic logic for optimal monitoring of security events with limited resources. J. Comput. Secur. 24, 6 (2016), 735–791.
- [12] Xavier Larriva-Novo, Mario Vega-Barbas, Victor A. Villagrá, Diego Rivera, Mario Sanz, and Manuel Álvarez-Campana. 2020. Dynamic risk management architecture based on heterogeneous data sources for enhancing the cyber situational awareness in organizations. In *Proceedings of the International Conference on Availability, Reliability and Security.* 1–9.
- [13] Greg Masters. June 9 2017. Crying Wolf: Combatting Cybersecurity Alert Fatigue. SC Magazine. Retreived from https://www.scmagazine.com/home/security-news/in-depth/crying-wolf-combatting-cybersecurity-alert-fatigue/.
- [14] Yuxin Meng and Lam-For Kwok. 2013. Enhancing false alarm reduction using voted ensemble selection in intrusion detection. Int. J. Comput. Intell. Syst. 6, 4 (2013), 626–638.
- [15] Tenda Okimoto, Naoto Ikegai, Katsumi Inoue, Hitoshi Okada, Tony Ribeiro, and Hiroshi Maruyama. 2013. Cyber security problem based on multi-objective distributed constraint optimization technique. In Proceedings of the 43rd Annual IEEE/IFIP Conference on Dependable Systems and Networks Workshop (DSN-W'13). IEEE, 1–7.
- [16] Ankit Shah, Rajesh Ganesan, and Sushil Jajodia. 2019. A methodology for ensuring fair allocation of CSOC effort for alert investigation. Int. J. Inf. Secur. 18, 2 (2019), 199–218.
- [17] Ankit Shah, Rajesh Ganesan, Sushil Jajodia, and Hasan Cam. 2018. Optimal assignment of sensors to analysts in a cybersecurity operations center. IEEE Syst. J. 13, 1 (2018), 1060–1071.
- [18] Ankit Shah, Rajesh Ganesan, Sushil Jajodia, Pierangela Samarati, and Hasan Cam. 2019. Adaptive alert management for balancing optimal performance among distributed CSOCs using reinforcement learning. *IEEE Trans. Parallel Distrib. Syst.* 31, 1 (2019), 16–33.
- [19] Yuan Wang, Yongjun Wang, Jing Liu, Zhijian Huang, and Peidai Xie. 2016. A survey of game theoretic methods for cyber security. In Proceedings of the IEEE 1st International Conference on Data Science in Cyberspace (DSC'16). IEEE, 631–636.
- [20] Mingtao Wu and Young Moon. 2019. Alert correlation for cyber-manufacturing intrusion detection. Proc. Manufact. 34 (2019), 820–831.

Received January 2021; revised September 2021; accepted September 2021

67:24