# Diversifying Database Activity Monitoring with Bandits

**Hagit Grushka-Cohen,[1] Ofer Biller,[2] Oded Sofer,[3] Lior Rokach,[1] Bracha Shapira,[1]**

[1]Dept. of Software and Information Systems Engineering,Ben-Gurion University of the Negev, Beer Sheva, Israel
[2]IBM Guardium Security Division and IBM Cyber Security Center of Excellence, Beer Sheva, Israel
[3]IBM Guardium Security Division, Tel-Aviv, Israel
hgrushka@post.bgu.ac.il, ofer.biller@il.ibm.com, odedso@il.ibm.com, liorrk@bgu.ac.il, bshapira@bgu.ac.il

## Abstract

Database activity monitoring (DAM) systems are commonly used by organizations to protect the organizational data, knowledge and intellectual properties. In order to protect organizations database DAM systems have two main roles, monitoring (documenting activity) and alerting to anomalous activity. Due to high-velocity streams and operating costs, such systems are restricted to examining only a sample of the activity. Current solutions use policies, manually crafted by experts, to decide which transactions to monitor and log. This limits the diversity of the data collected. Bandit algorithms, which use reward functions as the basis for optimization while adding diversity to the recommended set, have gained increased attention in recommendation systems for improving diversity.

In this work, we redefine the data sampling problem as a special case of the multi-armed bandit (MAB) problem and present a novel algorithm, which combines expert knowledge with random exploration. We analyze the effect of diversity on coverage and downstream event detection tasks using a simulated dataset. In doing so, we find that adding diversity to the sampling using the bandit-based approach works well for this task and maximizing population coverage without decreasing the quality in terms of issuing alerts about events.

## Introduction

Databases are at the heart of organizational IT infrastructure. For data security, privacy protection, and data leakage prevention, many organizations monitor database operations in real-time using database activity monitoring (DAM) systems. Such systems are widely used to help implement security policies, and detect attacks and data abuse. The multibillion dollar Google-Waymo vs. Uber case [1], in which the data collected (and more specifically, data activity logs) was used to document data theft, provides a good example of DAM system use and the value of such a system to an organization.

With information overload of hundreds of thousands of transactions per second, database activity monitoring systems cannot process and log all of the activity into log files.

Instead of logging all transactions, DAM systems use policies to decide which transactions to save, which are based on rules and users / activity groups defined by the SO during the system setup.

Policy changes require a lot of manual effort, resulting in a situation in which policies rarely change once they are defined. Furthermore, the system primarily monitors users that match the security officer's (SO) implicit or explicit preferences very closely, and this restricts them from discovering concept drifts (such as changes in users' roles).

This static approach may cause the "filter bubble" phenomenon in which SOs (the users of the DAM system) are trapped in a subspace of options that are too similar to the defined risk profile,thereby losing the ability to explore beyond what they already know. Such filter bubbles are a known issue for recommendation systems (Nguyen et al. 2014; Chen et al. 2019).

To address this challenge, we suggest incorporating the concept of diversity from recommendation systems (Matt et al. 2014) into logging policies. Unlike search engines or recommendations, sampling a more diverse group of users is not technically complicated as the user's transactions risk can be aggregated to a single score (Grushka-Cohen et al. 2016; Evina et al. 2019). However, logging capacity is constrained, and by focusing solely on diversity, undocumented malicious activity in the high risk group can be missed. A solution must provide a balance between the exploration of activities and users suspected to be of low risk with the exploitation of activities and users suspected to be of high risk.

The effect of sampling has been studied in the domain of anomaly detection systems of network traffic monitoring and DAM (Juba et al. 2015; Jadidi et al. 2016; Jadidi et al. 2015; Grushka-Cohen et al. 2017). These studies found that policy based sampling introduces bias to the down stream anomaly detection systems and suggested specific sampling strategies for avoiding bias.

Multi-armed bandits (MABs) are strategies used for policy setting based on sampling for decision-making (Agrawal and Goyal 2012; Auer, Cesa-Bianchi, and Fischer 2002). To the best of our knowledge, MABs haven't been used so far by an anomaly detection system to set data collection policy. Since auditing and anomaly detection system can only learn and model the data available to it, the policy determining which transactions are logged has great influ-

---

ence on the ability to detect anomalies, and which anomalies can be detected, and what actions would be possible to audit later. MABs are used for balancing exploration and exploitation in order to maximize a reward in many domains, including reinforcement learning and recommendation systems. In the stochastic multi-armed bandit problem, the goal is to maximize the reward from a slot machine with multiple arms. Various strategies for balancing exploration/exploitation have been studied (Agrawal and Goyal 2012).

In this work we suggest viewing the diversity problem for DAM system sampling strategies as a MAB problem, where the risk of the transactions logged is used as the reward function.

Unlike the classic MAB problem, in this special case, the risk distribution of a user is not static: it may change naturally when a user changes his/her role, or it can change due to hacking or malware, or when an employee has been compromised. Another difference is that in each round we are not sampling one user's activity but sampling multiple users to monitor, *i.e.* multiple arms are pulled in each round (the collector logs all of the transactions for a list of users).

We define a novel variant of the MAB problem to incorporate capacity, pulling multiple levers in each time frame, with a reward function which may change, named Budget-Constrained-Dynamic-MAB. We present a novel sampling strategy for sampling Budget-Constrained-Dynamic-MAB, a variant of the $\epsilon$–Greedy algorithm, named C–$\epsilon$–Greedy.

Using simulated data-sets created with (Grushka-Cohen et al. 2019), we assess the effect of diversity on sampling policies. We introduce two evaluation measures: (i) coverage of user activity, *i.e.*, how much do we know about each user in the population, and (ii) downstream effectiveness in identifying anomalies.

Using these measures, we compare a number of strong baselines and show that the proposed algorithm outperforms oracle-knowledge policy sampling and the Gibbs-prior sampling approach suggested in (Grushka-Cohen et al. 2019).

The novelty of our work lies in adding diversity as an important objective of monitoring and redefining the monitoring problem as a bandit problem, as well as in showing that a variant of a classic bandit sampling algorithm can improve existing approaches.

## Background and Motivation

### Database Monitoring

Monitoring database and file activity enables organizations to issue alerts when dangerous events have occurred, and database monitoring is implemented in many domains, including health, finance, and insurance. Security information and event management (SIEM) and DAM systems audit database activity and apply anomaly detection algorithms in an attempt to detect data misuse, data leakage, impersonation, and attacks on database systems. These systems are used by an organization's SO to identify suspicious activity (Kaplan, Sharma, and Weinberg 2011). Unlike recommendation systems where the users modeled are also the ones receiving the recommendation, in DAM settings the

aim model the database users' activity in order to make recommendations for the SO regarding which users should be monitored.

When monitoring only users suspected of preforming risky activity without exploring and diversifying the monitored user list, the SO can only learn about the "normal behaviour" of a small group of users without getting a sense of what is normal for the majority of users and activities.

### Transaction Risk

When an SO assigns risk to a transaction, various contextual information is used, such as time of day, user activity profile, location (IP address), the nature of the activity (*i.e.* is it permitted), data sensitivity, and the resulting data volume.

When a DAM system is installed in an organization these rules can be defined manually (by the SO) as a risk policy or learned by annotating risk scores on some representative transaction using a classifier such as CyberRank (Grushka-Cohen et al. 2016).

This allows for the aggregation of the overall user activity and the assignment of a single risk-score for user activity within a time-frame. Methods for aggregating this are out of scope of this work, but simple aggregations such as maximum or median risk during a particular hour, can be used.

### Increasing diversity when sampling for monitoring

The problem of data velocity and size when monitoring systems is not limited to the DAM domain, nor is the filter bubble effect. Studied conducted in the domain of network traffic flows such as (Jadidi et al. 2016; Mai et al. 2006) attempted to diversify the monitored data while maintaining the ability to detect malicious events. Mai *et al.* (Mai et al. 2006) found that for network flows the attempted sampling schemes weakened the anomaly detection ability to identify abrupt changes in volume - the sampling schemes degraded the performance examined in terms of success detection and false positive ratio. Kumar and Xu (Kumar and Xu 2006) suggested approximating the size of the traffic flow as a prior for the sampling probability, using count sketch for guiding the sampling method to estimate flow traffic measurement.

These studies found that sampling only high-risk packets is deleterious for anomaly detection and that sampling introduced bias to the recalls of the anomaly detection.

### DAM Sampling as a Multi Armed Bandit Problem

The problem of maximizing reward when sampling in the MAB problem has been extensively explored (May and Leslie 2011; Agrawal and Goyal 2012; Chen, Wang, and Yuan 2013). In multi-armed bandit problems, the optimal sampling (exploration / exploitation) strategy of playing against a slot machine with N arms each with an unknown reward distribution, is sought.

In the data-base monitoring domain, the goal is to find the optimal strategy for sampling users' database transactions, using the available resources in order to maximize risk monitoring. This does not fit the classical MAB setting as we are pulling multiple levers in every round. The scenario of making multiple actions at each round with a constraint on budget had been recently described by (Zhou and Tomlin 2018)

as Budget-Constrained MAB with Multiple Plays. However, in our case the reward distribution for each arm may change during the game and the rewards are not sparse as some reward is gained from every armed pulled (unlike vanilla MAB problems where most actions result in no reward). We name this MAB variant Budget-Constrained-Dynamic-MAB.

Once we redefine the security sampling problem as a MAB problem we can define the sampling algorithms in MAB terms of exploration and exploitation and use reward as the objective, aligned to the task of maximizing the collection of informative and useful logs.

## Problem Setting

We address the problem of diversifying the list of users who's transactions are being monitored based on the need of the system to monitor and capture all potential malicious activity, and the need of the SO to learn about the different users activity. Hence, we define the learning to diversify problem for database activity monitoring along with maximization of risk "reward" subject to the available storage and computing capacities. We now create a few notations to explain the problem. Let's consider a policy $p_t$ at a certain time iteration $t \in \{0, ...., \infty\}$ monitors $S_t$ subset of the system users $U = \{u_1, ..., u_n\}$ whos risk scores for time t are given by $\{r_1, ...r_n\}$. Our goal is to select the best policy $p_t$ (A subset of users from the n users) for each time-frame $t$ such that the result set caters to maximizing the monitored risk score subject to the given capacity $C$. Let us use indicator variables $x_{jt} = \{0, 1\}$ denotes if the user $u_j$ is selected by policy $p_t$ to be monitored during time-frame t, and $z_{jt} = \{0, 1\}$ denotes if the user $u_j$ is selected by the oracle policy $o_t$ to be monitored during time-frame t . The reward generated from the policy $p_t$ can be computed as $\rho_t = \sum_{j=1}^{j=n} x_{jt}r_{jt}$. The reward proportion $r_t$ for time t can be computed by

$$R_t = \frac{\rho_t}{\sum_{j=1}^{j=n} z_{jt}r_{jt}}$$

The total reward can then be given by $R_T = \sum_{t=1}^{t=T} R_t$. Given the above set up learning to diversify problem in Database activity monitoring can be mapped to maximizing $R_T$ for a given capacity $C$ without degrading the Anomalies found.

Each users activity is either fully monitored or fully discard.The anomaly detection aims to detect advanced persistent threat (APT) (Tankard 2011). a single risk score for the users' time-frame activity lines up with these goals.

## Evaluation

In this study, we use the following metrics to evaluate the performance of sampling algorithms from the perspective of risk detection, coverage, and malicious events detected:

- **Reward** The main objective of the sampling strategy is to maximize the monitored risk . The oracle strategy, detects the maximal risk $\rho_{opt} = \sum_{j=1}^{j=n} z_{jt}r_{jt}$ for a given capacity at time t. The reward for a time-frame is the proportion of the detected risk out of the oracle strategy detected risk: $R_t = \frac{\rho_t}{\rho_{oracle}}$

- **Coverage** Anomaly detection systems rely on historical data in order to detect anomalies. Therefore acquiring risk profiles for as many users as fast as possible is highly important, also for avoiding the cold start problem. If we never sample the users suspected to be a low risk users the anomaly detection wouldn't be able to detect a change in their activity. Moreover, the SO perception of the user's risk profile may be inaccurate or just missing. We measure how long does it take the system to collect data regarding the entire population. Therefore, we measure the percent of the population for which the sampling algorithm logged at least one time frame of activity. Let us define $users\_coverage_{t_i} = \{U_t \mid t \in \{1, ..., t_i\}\}$ where $U_t$ is a list of all users monitored during time $t$

- **Recall of anomalies** Downstream of the logged data, anomaly detection is applied to user activity. Previous work (Chandola, Banerjee, and Kumar 2009) found that sampling adds bias to the results of the anomaly detection. We measure the recall of an identical anomaly detection algorithm when applied to the different data sets produced by each sampling strategies. We normalize the recall to the recall discovered by applying the anomaly detection over the full data set (no sampling).

## Sampling algorithms

We compare three algorithms for the problem. The first two baseline algorithms proposed in (Grushka-Cohen et al. 2019). We introduce a novel algorithm named C–$\epsilon$–Greedy based on $\epsilon$-greedy algorithm for the Budget-Constrained-Dynamic-MAB problem.

### SO-policy

In each round, the C users with the highest known risk are sampled (these remain constant as no exploration is performed). Initialization is explained in the Experimental settings section.

### Gibbs-by-risk

In the Gibbs-exploration strategy, suggested by (Grushka-Cohen et al. 2019), the user probability to be sampled is proportional to the risk estimate.

---

**Algorithm 1** Gibbs by risk

---

**Data:** $U = \{u_1, .., u_n\}$, $R_{t-1} = \{r_{1,t-1}, .., r_{n,t-1}\}$, $C > 0$
**Result:** list of users for time $t$

1   monitor_users $= \{\}$
2   **for** *u in $\{u_1,...,u_n\}$* **do**
3     $p_u = r_{u_i} / \sum_{i=1}^{N} r_{u_i}$
4   **end**
5   **while** *len(monitor_users) $\leq C$* **do**
6     draw $u_i$ with probability $p_i \sim (p_1, ...p_n)$
7     **if** $u_i \notin$ *monitor_users* **then**
8       monitor_users.append($u_i$)
9     **end**
10   **end**

---

Sample randomly with a user's risk profile as the prior for choosing it. The user's probability to be sampled is determined by the maximal user's risk found in the past (using a sliding window of k time frames) or the initialization value.

## C-$\epsilon$-greedy

The $\epsilon$–greedy strategy for MAB has been shown to outperform other algorithms on most settings (Kuleshov and Precup 2014). Thompson sampling, another top performing algorithm for MAB, is not applicable to the non-discrete nature of Budget-Constrained-Dynamic-MAB. In our setting we present C–$\epsilon$–Greedy strategy, where in each step the $\epsilon \times C$ (capacity) users with the highest risk are sampled, and $(1 - \epsilon) \times C$ random users are sampled as well for exploration.

---

**Algorithm 2** C-$\epsilon$-greedy

**Data:** $U = \{u_1, .., u_n\}, R_{i,t} = \{r_{1,1}, .., r_{n,n}\}, C > 0, k > 0$

**Result:** list of users for time t

11 monitor_users ={}
12 **for** $u_i \in U$ **do**
13 $\quad aggrisk_{u_i} = \sum_{t=T-k}^{T} r_{u_i,t}$
14 $\quad p_{u_i} = aggrisk_{u_i}/k$
15 **end**
16 sorted_users $\leftarrow sort(((u_1, p_1), .., (u_n, p_n)), desc)$
17 $i = 0$
18 **while** $i < \epsilon \times C$ **do**
19 $\quad$ monitor_users.append(sorted_users[i][0])
20 $\quad i+ = 1$
21 **end**
22 $i = 0$
23 **while** $i < (1 - \epsilon) \times C$ **do**
24 $\quad$ draw j with probability $U \sim (1, .., n)$
25 $\quad$ **if** *sorted_users[j][0]* $\notin$ *monitor_users* **then**
26 $\quad\quad$ monitor_users.append(sorted_users[j][0]) $i+ = 1$
27 $\quad$ **end**
28 **end**

---

The SO-Policy can be viewed as a special case of the C–$\epsilon$–Greedy Strategy where the exploitation is set to 100%.

### Random sampling

A baseline where users are sampled completely randomly at every time frame. This can be viewed as a special case of the C–$\epsilon$–Greedy Strategy where the exploration is set to 100%.

## Simulating DAM monitoring data

Using data collected from an operating DAM system is irrelevant in our case mainly because it contains bias introduced by the process of the data collection; Due to the high velocity nature of database activity, it is impossible to log all the activity, therefore the data collected using a sampling policy, hence contains bias. Other reasons we worked with simulated data: (1) The logged data may contain organizational Intellectual property (IP) and sensitive data and as such it is

hard to convince commercial organization to share it. (2) Security events are not tagged in the logged data, therefore for each suspected activity a consultation with SO is needed. Verifying that the activity is malicious requires manual investigation of the activity and it's context making it impractical for research purposes. Using simulated data allows us to evaluate the effect of the sampling strategy using all the evaluation metrics described above.

Collecting meaningful and unbiased data sets for studying effects of diversity and comparing sampling strategies is also a problem in other domains (Chen et al. 2019; Lupu, Durand, and Precup 2019) such as search and recommendations, using simulated data instead allows for flexible and robust experimentation (Goswami, Zhai, and Mohapatra 2019; Chen et al. 2019; Germano, Gómez, and Le Mens 2019; Voelkl et al. 2018).

DAM systems log various features for each user transaction, using the data as is makes evaluating log coverage and anomaly detection more complex as it is not a trivial task. To reduce the complexity of the features for comparison and evaluation, working with aggregated data is useful. Previous work such as (Grushka-Cohen et al. 2016; Evina et al. 2019) leveraged SO knowledge to aggregate database activity into a single risk score. (Grushka-Cohen et al. 2019) suggested a simulation package made of low complexity data where the user activity for a time frame is represented by a single aggregated risk score. Using simulated data also helps us with avoiding the bias introduced to the data set during the collecting stage and the ability to introduce anomalies in a controlled fashion.

The simulation is designed to mimic the properties of real users whose behavior is not random: users exhibit power log distribution of risk profile, few users produce most of the risky transactions, and transactions per user over time have a trend (see (Grushka-Cohen et al. 2019)). Before implementing our experiment we consulted a security expert to confirm that the simulated data matches the expected behaviour as observed in real data-bases activity.
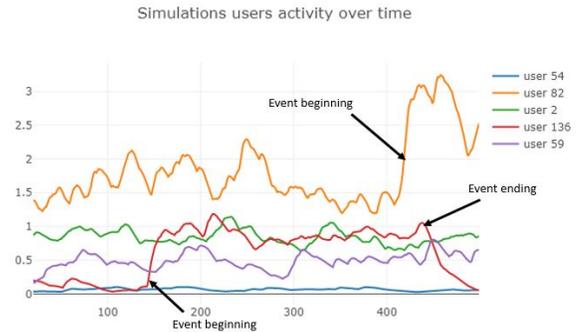


Figure 1: Risk distribution. Risk sampled per user over time produced by the simulation, two of the users have an anomaly introduced

## Experimental settings

We used the simulation to produce 10 data sets using different random seeds and report on the average results. Each data set simulates activity profiles of 200 users for 3,000 time frames. The capacity (C) was set to 10% of the numbers of users.

## Initialization

The setup of a DAM system relies on the SO's acquaintance with the databases users as well as his/her familiarity with the domain potential risks. During the setup process the SO defines a monitoring policy made of rules and user groups which represent the users and activities suspected of being the main threats to the organization's data. To examine the effect of that knowledge over the quality of the monitoring and anomaly detection and to assess whether the results of the SO's manual efforts, can be leveraged as a prior for the cold start scenario we compare two settings:

(1) A setting in which there is *oracle initial knowledge* - This setting assumes that the SO, defining the policy, has perfect knowledge. To mimic this knowledge we sample the risk for each user at time $t_0$ and use it as the risk prior. (2) A setting in which there is a noisy oracle - this settings assumes that the SO has imperfect knowledge about the users' activity and the risk they present to the organization. In such cases, we use the risk from time $t_0$ mixed with noise from another randomly generated user (risk distribution) as a risk prior for each user.

## Introducing security events

The anomaly detection component of DAM systems relies on modeling users' activity and detecting activities that are incompatible with that distribution. To evaluate the effect of the sampling strategy on the anomaly detection component, we simulated security events. A security event, in which a user has been compromised or abused his/her permissions, is continuous (lasts more than a single time frame) and characterized by a change in the user's risk distribution.

To simulate such an event we change the user's risk distribution, sampling a new risk distribution randomly for a period length $l$ which is sampled uniformly between 200 and 300. In each time frame, a user has a probability of experiencing such an event (set at 0.001 for these experiments).

Figure 1, presents the activity of five different users from the users generated by the simulator. Each line indicates the users' hourly activity risk score. The graph presents one high-risk profiled user (user 82), two medium-risk profiled users (users 2 and 59), and two low-risk profiled users (users 136 and 54). User 82 experienced a security event beginning at the 414 time frame. User 136 experienced a security event between the 143 to 441 time frames. In the figure it can be seen that user 143 started of as a low risk user, this user's activity profile changed to a medium-risk profile during the security event, while user 82 had a high-risk profile to begin with and the security event exacerbated this, resulting in an even higher-risk profile for this user.
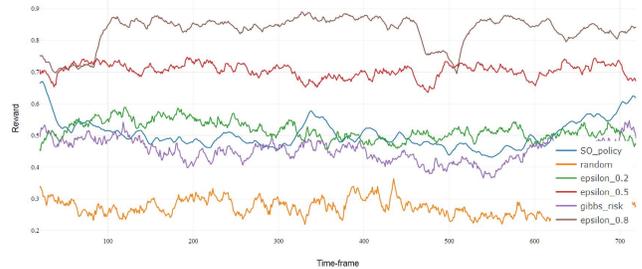


Figure 2: Risk Reward per Time Frame

## Controlling diversity

To evaluate the effect of different amounts of introduced diversity we tweak the $\epsilon$ part of the C–$\epsilon$–Greedy algorithm. When $\epsilon$ is 1 only the highest risk users are sampled and no exploration is happening, this is essentially the SO-policy baseline. Setting $\epsilon$ to 0 means the system is in pure exploration mode, sampling users at random. 3 other settings are evaluated: setting $\epsilon$ to 0.2 (sample 80% of capacity at random), setting $\epsilon$ to 0.5 (sample 50% of capacity at random and 50% from the highest risk users) and setting $\epsilon$ to 0.8 (sample only 20% of capacity at random),

## Results

### Performance

Table 1 presents the performance for the reward and recall metrics when using the various sampling strategies. As can be seen, the C–$\epsilon$–Greedy sampling strategy, with exploration set at 20% or 50% ($\epsilon$ set at 80% or 50%) out-performed all of the other strategies. C–$\epsilon$–Greedy, with exploration of 20% ($\epsilon$ set at 80%) performed the best in terms of reward, with reward of 86%. when initialized With a perfect SO knowledge (assuming the SO has the right estimation of all the users activity during setup) the SO-Policy of sampling only the riskiest users achieved a 67% of the optimal reward, however when the initialization is noisy (imperfect expert knowledge at system setup) the result degraded to 51% of the reward while the C–$\epsilon$–Greedy methods did not degrade significantly. Figure 2, which describes the reward achieved in each time frame, shows that the C–$\epsilon$–Greedy sampling strategy both with the 80% / %20 exploitation / exploration ratio and 50% / 50% exploitation / exploration ratio, outperforms all other strategies, while random sampling constantly results in the poorest performance. In the figure, we see that C–$\epsilon$–Greedy 80% (using 20% random exploration) takes about 100 time frames to explore, discovering the most risky users and then exploiting gaining  20% better reward than C–$\epsilon$–Greedy 50%. At time frame 460 we see another event where a change in the user population, a non-risky user affected by an anomaly that made it emit high risk transactions, cause both C–$\epsilon$–Greedy 80% and 50% to experience a drop in reward, the reward goes back up once the offending user is discovered and added to the exploitation list.

In terms of anomalies recall, the SO-knowledge policy detected only 8.5% of the anomalies, not surprising as only

Table 1: Reward and Recall for sampling strategies

| Strategy | Reward (detected risk, Oracle SO) | Reward (with Noisy-Oracle SO) | Recall (detected anomalies) |
|---|---|---|---|
| SO policy | 0.677 | 0.514 | 0.084 |
| Random | 0.264 | 0.264 | **0.878** |
| Gibbs by Risk | 0.514 | 0.39 | 0.56 |
| C–$\epsilon$–Greedy 0.2 | 0.479 | 0.48 | 0.834 |
| C–$\epsilon$–Greedy 0.5 | 0.703 | 0.702 | 0.779 |
| C–$\epsilon$–Greedy 0.8 | **0.864** | **0.863** | 0.67 |

10% of the users are ever monitored. The highest recall was found for random sampling detected 88% of the anomalies. The C–$\epsilon$–Greedy sampling strategy detected between 67% to 83% of the anomalies. Gibbs-by-risk sampling got significantly less recall with 56%.

## Coverage experiment

In figure 3 we compare the coverage of the different algorithms. On the x axis we shows the time and the Y axis shows the corresponding coverage.

We marked a user as "covered" when the method gathered two samples of the user.

The SO-policy always samples the same group of users, therefore it is constant from the beginning as no other users are explored. With Gibbs-by-risk method there's a probability for any user to be sampled in proportion to their risk. We observe that users who exhibit low risk in initialization are not likely to be sampled which slows down the exploration of all users (not all users were sampled even after 3,000 time frames).

The strategies preferring completely random exploration had faster knowledge acquisition, depending on the proportion of the sample capacity devoted to exploration. When sampling with a 100% exploration rate (completely random), 80% exploration (epsilon_0.2), or 50% exploration the strategy had gathered two samples for most users by the 100th time frame. Dedicating 20% of the capacity to exploration (epsilon_0.8) requires ∼ 250 time frames for the same amount of knowledge acquisition.
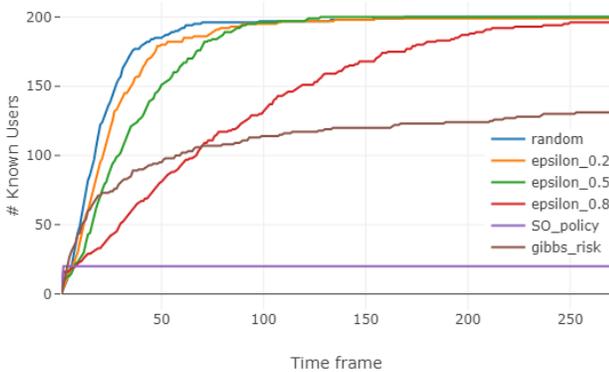


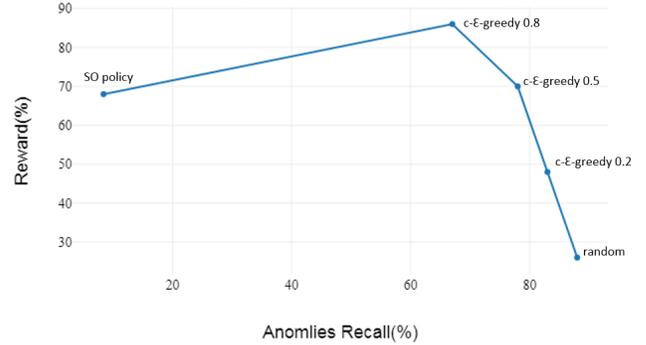Figure 3: Knowledge acquisition rate



Figure 4: Reward vs Recall as a function of exploration rate

## Anomaly detection recall vs. reward

Figure 4 shows that anomalies recall improves as we increase the exploration rate. However, the highest exploration rate yields low reward.

## Conclusions

In this paper, we explored the effects of diversity when sampling activity for monitoring. We formulated the problem of sampling transactions in the security domain as a MAB problem, introducing Budget-Constrained-Dynamic-MAB, a multi-armed bandit where in each time-frame we sample a number of arms defined by the capacity of the system. As a user's risk profile may change when they switch position or become compromised, the reward probability may change during the user's life cycle. We extended the $\epsilon$–Greedy algorithm, one of the best performing MAB approaches (Chen, Wang, and Yuan 2013), to Budget-Constrained-Dynamic-MAB. Our variant, C–$\epsilon$–Greedy splits the sampling capacity at each turn into exploration part and exploitation part.

Using simulated data sets of user activity we showed that ensuring diversity when sampling database activity is important for understanding user behavior (coverage) and enhances downstream anomaly detection without adverse effects from allocating some capacity to diversity.

The C–$\epsilon$–Greedy strategy that used 20% of the capacity for exploration was able to beat a strong oracle-policy baseline on collecting risky activity, while achieving good coverage of all the population. The baseline method of a constant policy did achieve fine results on reward, collecting logs of risky activity, but was sensitive to initialization (depending on the SO familiarity with all the users). C–$\epsilon$–Greedy having better coverage allowed it to identify users becoming riskier,

achieving better reward overall and was immune to bad initialization.

We found that recall is driven by the exploration part of the strategy, and that completely random exploration had the best recall but worst reward (most of the logging capacity was spent on non-risky activity).

Therefore, we advise using methods for sampling that allow the SO to dedicate some capacity for exploration and diversity. We find that the The C–$\epsilon$–Greedy strategy is easy to implement, the $\epsilon$ is an easy knob to turn for enhancing exploration when needed, it supports initialization with the SO knowledge avoiding cold start and keeping humans in the loop while maximizing the audit trail for all users and supporting robust anomaly detection downstream.

# References

[Agrawal and Goyal 2012] Agrawal, S., and Goyal, N. 2012. Analysis of thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, 39–1.

[Auer, Cesa-Bianchi, and Fischer 2002] Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-time analysis of the multiarmed bandit problem. *Machine learning* 47(2-3):235–256.

[Chandola, Banerjee, and Kumar 2009] Chandola, V.; Banerjee, A.; and Kumar, V. 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41(3):15.

[Chen et al. 2019] Chen, L.; Yang, Y.; Wang, N.; Yang, K.; and Yuan, Q. 2019. How serendipity improves user satisfaction with recommendations? a large-scale user evaluation. In *The World Wide Web Conference*, 240–250. ACM.

[Chen, Wang, and Yuan 2013] Chen, W.; Wang, Y.; and Yuan, Y. 2013. Combinatorial multi-armed bandit: General framework and applications. In *International Conference on Machine Learning*, 151–159.

[Evina et al. 2019] Evina, P. A.; AYACHI, F. L.; JAIDI, F.; and BOUHOULA, A. 2019. Enforcing a risk assessment approach in access control policies management: Analysis, correlation study and model enhancement. In *2019 15th International Wireless Communications & Mobile Computing Conference (IWCMC)*, 1866–1871. IEEE.

[Germano, Gómez, and Le Mens 2019] Germano, F.; Gómez, V.; and Le Mens, G. 2019. The few-get-richer: a surprising consequence of popularity-based rankings? In *The World Wide Web Conference*, 2764–2770. ACM.

[Goswami, Zhai, and Mohapatra 2019] Goswami, A.; Zhai, C.; and Mohapatra, P. 2019. Learning to diversify for e-commerce search with multi-armed bandit.

[Grushka-Cohen et al. 2016] Grushka-Cohen, H.; Sofer, O.; Biller, O.; Shapira, B.; and Rokach, L. 2016. Cyberrank: Knowledge elicitation for risk assessment of database security. 2009–2012. ACM.

[Grushka-Cohen et al. 2017] Grushka-Cohen, H.; Sofer, O.; Biller, O.; Dymshits, M.; Rokach, L.; and Shapira, B. 2017. Sampling high throughput data for anomaly detection of data-base activity. *arXiv preprint arXiv:1708.04278*.

[Grushka-Cohen et al. 2019] Grushka-Cohen, H.; Biller, O.; Sofer, O.; Rokach, L.; and Shapira, B. 2019. Simulating user activity for assessing effect of sampling on db activity monitoring anomaly detection. In *Policy-Based Autonomic Data Governance*. Springer. 82–90.

[Jadidi et al. 2015] Jadidi, Z.; Muthukkumarasamy, V.; Sithirasenan, E.; and Singh, K. 2015. Performance of flow-based anomaly detection in sampled traffic. *Journal of Networks* 10(9):512.

[Jadidi et al. 2016] Jadidi, Z.; Muthukkumarasamy, V.; Sithirasenan, E.; and Singh, K. 2016. Intelligent sampling using an optimized neural network. *Journal of Networks* 11(01):16–27.

[Juba et al. 2015] Juba, B.; Musco, C.; Long, F.; Sidiroglou-Douskos, S.; and Rinard, M. C. 2015. Principled sampling for anomaly detection. In *NDSS*.

[Kaplan, Sharma, and Weinberg 2011] Kaplan, J.; Sharma, S.; and Weinberg, A. 2011. Meeting the cybersecurity challenge. *Digit. McKinsey Google Scholar*.

[Kuleshov and Precup 2014] Kuleshov, V., and Precup, D. 2014. Algorithms for multi-armed bandit problems. *arXiv preprint arXiv:1402.6028*.

[Kumar and Xu 2006] Kumar, A., and Xu, J. J. 2006. Sketch guided sampling-using on-line estimates of flow size for adaptive data collection. In *Infocom*.

[Lupu, Durand, and Precup 2019] Lupu, A.; Durand, A.; and Precup, D. 2019. Leveraging observations in bandits: Between risks and benefits.

[Mai et al. 2006] Mai, J.; Chuah, C.-N.; Sridharan, A.; Ye, T.; and Zang, H. 2006. Is sampled data sufficient for anomaly detection? In *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*, 165–176. ACM.

[Matt et al. 2014] Matt, C.; Benlian, A.; Hess, T.; and Weiß, C. 2014. Escaping from the filter bubble? the effects of novelty and serendipity on users' evaluations of online recommendations.

[May and Leslie 2011] May, B. C., and Leslie, D. S. 2011. Simulation studies in optimistic bayesian sampling in contextual-bandit problems. *Statistics Group, Department of Mathematics, University of Bristol* 11:02.

[Nguyen et al. 2014] Nguyen, T. T.; Hui, P.-M.; Harper, F. M.; Terveen, L.; and Konstan, J. A. 2014. Exploring the filter bubble: the effect of using recommender systems on content diversity. In *Proceedings of the 23rd international conference on World wide web*, 677–686. ACM.

[Tankard 2011] Tankard, C. 2011. Advanced persistent threats and how to monitor and deter them. *Network security* 2011(8):16–19.

[Voelkl et al. 2018] Voelkl, B.; Vogt, L.; Sena, E. S.; and Würbel, H. 2018. Reproducibility of preclinical animal research improves with heterogeneity of study samples. *PLoS biology* 16(2):e2003693.

[Zhou and Tomlin 2018] Zhou, D. P., and Tomlin, C. J. 2018. Budget-constrained multi-armed bandits with mul-

tiple plays. In *Thirty-Second AAAI Conference on Artificial Intelligence*.