# A Scalable Algorithm for Anomaly Detection via Learning-Based Controlled Sensing

Geethu Joseph, M. Cenk Gursoy, and Pramod K. Varshney, *Life Fellow, IEEE*

*Dept. of Electrical Engineering and Computer Science*

*Syracuse University*

Syracuse, NY 13244, USA

Emails:{gjoseph,mcgursoy,varshney}@syr.edu.

*Abstract*—We address the problem of sequentially selecting and observing processes from a given set to find the anomalies among them. The decision maker observes one process at a time and obtains a noisy binary indicator of whether or not the corresponding process is anomalous. In this setting, we develop an anomaly detection algorithm that chooses the process to be observed at a given time instant, decides when to stop taking observations, and makes a decision regarding the anomalous processes. The objective of the detection algorithm is to arrive at a decision with an accuracy exceeding a desired value while minimizing the delay in decision making. Our algorithm relies on a Markov decision process defined using the marginal probability of each process being normal or anomalous, conditioned on the observations. We implement the detection algorithm using the deep actor-critic reinforcement learning framework. Unlike prior work on this topic that has exponential complexity in the number of processes, our algorithm has computational and memory requirements that are both polynomial in the number of processes. We demonstrate the efficacy of our algorithm using numerical experiments by comparing it with the state-of-the-art methods.

*Index Terms*—Active hypothesis testing, anomaly detection, deep learning, reinforcement learning, actor-critic algorithm, quickest state estimation, sequential decision-making, sequential sensing.

## I. INTRODUCTION

We consider the problem of observing a given set of processes to detect the anomalies among them via controlled sensing. Here, the decision maker does not observe all the processes at each time instant, but sequentially selects and observes one process at a time. The sequential control of the observation process is referred to as controlled sensing. The challenge here is to devise a selection policy to sequentially choose the processes to be observed so that the decision is accurate and fast. This problem arises, for instance, in sensor networks used for remote health monitoring, structural health monitoring, etc [1], [2]. Such systems are equipped with different types of sensors to monitor different functionalities (or processes) of the system. The sensors send their measurements to a common decision maker that identifies any potential system malfunction. These sensor measurements can be noisy due to faulty hardware or unreliable communication links. Therefore, to ensure the accuracy of the decision, we employ a sequential process selection strategy that observes the set of processes one at a time over multiple time instants before the final decision is made. Further, the different processes can be

dependent on each other, and therefore, observing one process also gives information about other dependent processes. Our goal is to derive a selection policy that accurately identifies the anomalous processes with minimum delay by exploiting the underlying statistical dependence among the processes.

A popular approach for solving the anomaly detection problem is to use the active hypothesis testing framework [3], [4]. Here, the decision maker defines a hypothesis corresponding to each of the possible states of the processes and computes the posterior probabilities over the hypothesis set using the observations. The decision maker continues to collect observations until the probability corresponding to one of the hypotheses exceeds the desired confidence level. This framework of active hypothesis testing was introduced by Chernoff in [5], and it was followed by several other studies in the literature [6]–[9]. Recently, some researchers have combined the active hypothesis testing framework with deep learning algorithms to design data-driven anomaly detection algorithms [3], [4], [10], [11]. These algorithms learn from a training dataset and come with an added advantage of adaptability to the underlying statistical dependence among the processes. The state-of-the-art algorithms in this direction employ reinforcement learning (RL) algorithms such as Q-learning [10] and actor-critic [3], [4], and the active inference framework [11]. However, the major drawback of this solution strategy is the heavy computational burden that arises due to the large number of hypotheses. Since each process can either be normal or anomalous, the number of hypotheses increases exponentially with the number of processes. Therefore, in this paper, we attempt to devise a learning-based controlled sensing framework for anomaly detection with polynomial complexity in the number of processes.

The specific contributions of the paper are as follows: we first reformulate the problem of anomaly detection in terms of the marginal (not joint) probability of each process being normal or anomalous, conditioned on the observations. Consequently, the number of posterior probabilities computed by the algorithm at every time instant is linear in the number of processes. Based on these marginal posterior probabilities, we define the notion of a confidence level that is proportional to the decision accuracy, and a reward function that monotonically increases with the decision accuracy and decreases with the duration of the observation acquisition phase. These

definitions allow us to reformulate the anomaly detection problem as a long-term average reward maximization of a Markov decision process (MDP). This problem is solved using a policy gradient RL algorithm called the actor-critic method, and the algorithm is implemented using deep neural networks. Using numerical results, we show that our algorithm is able to learn and adapt to the statistical dependence among the processes. Further, the polynomial complexity of the algorithms makes it scalable, and hence, practically more useful.

## II. ANOMALY DETECTION PROBLEM

We consider a set of $N$ processes where the state of each process is a binary random variable. The process state vector is denoted by $\boldsymbol{s} \in \{0,1\}^N$ whose $i^{\text{th}}$ entry being 0 and 1 indicates that the $i^{\text{th}}$ process is in the normal state and the anomalous state, respectively. We aim to detect the anomalous processes, which is equivalent to estimating the random binary vector $\boldsymbol{s}$.

We estimate the process state vector $\boldsymbol{s}$ by selecting and observing one process at every time instant, and obtaining a state estimate of the corresponding process which has a finite probability of being erroneous. Let the process observed at time $k$ be $a(k) \in \{1, 2, \ldots, N\}$ and the corresponding observation be $y_{a(k)}(k) \in \{0,1\}$. The uncertainty in the observation is modeled using the following probabilistic model:

$$y_{a(k)}(k) = \begin{cases} s_{a(k)} & \text{with probability } 1-p, \\ 1 - s_{a(k)} & \text{with probability } p, \end{cases} \quad (1)$$

where $p \in [0,1]$ is called the flipping probability. Further, we assume that conditioned on the value of $\boldsymbol{s}$, the observations obtained across different time instants are jointly (conditionally) independent, i.e., for any $k$,

$$\mathbb{P}\left[\{y_i(l), i = 1, 2, \ldots, N\}_{l=1}^k \Big| \boldsymbol{s}\right] = \prod_{i=1}^N \prod_{l=1}^k \mathbb{P}\left[y_i(l)|s_i\right]. \quad (2)$$

Therefore, the $i^{\text{th}}$ process $\{\boldsymbol{y}_i(k) \in \{0,1\}\}_{k=1}^{\infty}$ is a sequence of independent and identically distributed (i.i.d.) binary random variables parameterized by $\boldsymbol{s}_i \in \{0,1\}$.

After each observation arrives, the decision maker computes an estimate of $\boldsymbol{s}$ along with the confidence in the estimate. The decision maker continues to observe the processes until the confidence exceeds the desired level denoted by $\pi_{\text{upper}} \in (0,1)$. Therefore, we have two interrelated tasks: one, to develop an algorithm to estimate the process state vector and the associated confidence in the estimate; and two, to derive a policy that decides the process to be observed at each time instant and the criterion to stop collecting observations. We seek the estimation algorithm and the policy that jointly minimize the stopping time $K$ while maximizing the accuracy level. Here, the stopping time refers to the time instant at which the observation acquisition phase ends. We next present our estimation algorithm and policy design.

## III. ESTIMATION ALGORITHM

In this section, we derive an algorithm to estimate the process state vector from the observations. We note that the observations depend on the selection policy, and the policy design, in turn, depends on the estimation algorithm. Therefore, we first present the estimation algorithm and then derive a selection policy based on the estimation objectives in the next section.

To estimate the process state vector, we first compute the belief vector $\boldsymbol{\sigma}(k) \in [0,1]^N$ at time $k$ whose $i^{\text{th}}$ entry $\sigma_i(k)$ is the posterior probability that the $i^{\text{th}}$ process is normal ($s_i = 0$). Therefore, the probability that the $i^{\text{th}}$ process is anomalous ($s_i = 1$) is $1 - \sigma_i(k)$. As each observation arrives, we recursively update the belief vector as follows.

$$\sigma_i(k) = \mathbb{P}\left[s_i = 0 \Big| \{y_{a(l)}(l)\}_{l=1}^k\right]$$
$$= \frac{\mathbb{P}\left[\{y_{a(l)}(l)\}_{l=1}^k \Big| s_i = 0\right] \mathbb{P}[s_i = 0]}{\mathbb{P}\left[\{y_{a(l)}(l)\}_{l=1}^k\right]}. \quad (3)$$

Here, we approximate the joint probability distribution by assuming that the observation $y_{a(k)}(k)$ is independent of the past observations $\{y_{a(l)}(l)\}_{l=1}^{k-1}$ conditioned on the process state $s_i$:

$$\mathbb{P}\left[\{y_{a(l)}(l)\}_{l=1}^k \Big| s_i = 0\right] \mathbb{P}[s_i = 0]$$
$$\approx \mathbb{P}\left[\{y_{a(l)}(l)\}_{l=1}^{k-1} \Big| s_i = 0\right] \mathbb{P}\left[y_{a(k)}(k)\big|s_i = 0\right] \mathbb{P}[s_i = 0]$$
$$= \sigma_i(k-1)\mathbb{P}\left[\{y_{a(l)}(l)\}_{l=1}^{k-1}\right] \mathbb{P}\left[y_{a(k)}(k)\big|s_i = 0\right]. \quad (4)$$

From (2), the observation $y_{a(k)}(k)$ is independent of all other observations, conditioned on the value of $s_{a(k)}$. Therefore, the approximation is exact when $s_{a(k)}$ is a deterministic function of $s_i$. Some examples of such cases are $\mathbb{P}\left[s_{a(k)} = s_i\right] = 1$, and $\mathbb{P}\left[s_{a(k)} = 1 - s_i\right] = 1$.

Substituting (4) into (3), we obtain

$$\sigma_i(k) = \frac{\sigma_i(k-1)\mathbb{P}\left[y_{a(k)}(k)\big|s_i = 0\right]}{\Sigma_i(k)}, \quad (5)$$

Here, following the approximation in (4), the normalization constant is

$$\Sigma_i(k) = \sigma_i(k-1)\mathbb{P}\left[y_{a(k)}(k)\big|s_i = 0\right]$$
$$+ (1 - \sigma_i(k-1))\mathbb{P}\left[y_{a(k)}(k)\big|s_i = 1\right]. \quad (6)$$

Further, the conditional probability $\mathbb{P}\left[y_{a(k)}(k)\big|s_i = s\right]$ for $s = 0, 1$ is given by

$$\mathbb{P}\left[y_{a(k)}(k)\big|s_i = s\right]$$
$$= \sum_{s'=0,1} \mathbb{P}\left[y_{a(k)}(k)\big|s_{a(k)} = s'\right] \mathbb{P}\left[s_{a(k)} = s'\big|s_i = s\right]$$
$$= \sum_{s'=0,1} \left[p^{|s' - y_{a(k)}(k)|}(1-p)^{|1-s'-y_{a(k)}(k)|}\right.$$
$$\left. \times \mathbb{P}\left[s_{a(k)} = s'\big|s_i = s\right]\right], \quad (7)$$

which follows from (1). Here, the term $\mathbb{P}\left[s_j = s' | s_i = s\right]$ can be easily estimated from the training data[1] for every pair $(i, j)$. Hence, (5), (6), and (7) give the recursive update of $\boldsymbol{\sigma}(k)$.

We note that when $s_{a(k)}$ and $s_i$ are independent processes,

$$\mathbb{P}\left[y_{a(k)}(k) | s_i = s\right] = \mathbb{P}\left[y_{a(k)}(k)\right] \quad s = 0, 1. \quad (8)$$

Consequently, (5) reduces to $\sigma_i(k) = \sigma_i(k-1)$. This update is intuitive since an observation from process $s_{a(k)}$ does not change the probabilities associated with an independent process $s_i$. In other words, the recursive relation is exact when $s_i$ and $s_{a(k)}$ are either independent or $s_{a(k)}$ can be exactly determined from $s_i$. We discuss this point in detail in Sec. VI.

Once $\boldsymbol{\sigma}(k)$ is obtained, the computation of the process state vector estimate denoted by $\hat{\boldsymbol{s}}(k)$ is straightforward:

$$\hat{s}_i = \begin{cases} 0 & \text{if } \sigma_i(k) \geq 1 - \sigma_i(k) \\ 1 & \text{if } \sigma_i(k) < 1 - \sigma_i(k). \end{cases} \quad (9)$$

Hence, the derivation of the estimation algorithm is complete. We next discuss the design of the selection policy.

## IV. Selection Policy

The design of the selection policy is a sequential decision making problem, and therefore, this problem can be formulated using the mathematical framework of an MDP. This formulation allows us to obtain the selection policy via reward maximization of the MDP using RL algorithms. In the following subsections, we define the MDP framework and describe the RL algorithm using the deep actor-critic method.

### A. Markov Decision Process

An MDP has four components: state space, action space, state transition probabilities, reward function. In our case, these components are defined as follows:

- *MDP state:* Our estimation algorithm is based on the belief vector $\boldsymbol{\sigma}(k)$ that changes with time after each observation arrives. Therefore, we define $\boldsymbol{\sigma}(k) \in [0, 1]^N$ as the state of the MDP at time $k$. We note that the MDP state vector $\boldsymbol{\sigma}(k)$ is different from the process state vector $\boldsymbol{s}$.
- *Action:* The state of MDP depends on the observation which in turn depends on the process selected by the policy. Thus, the action taken by the decision maker at time instant $k$ is the selected process $a(k) \in \{1, 2, \dots, N\}$.
- *MDP State Transition:* For our problem, the MDP state $\boldsymbol{\sigma}(k)$ at time $k$ is a deterministic function of the previous MDP state $\boldsymbol{\sigma}(k-1)$, the action $a(k)$, and the observation $y_{a(k)}(k)$. Therefore, the MDP state transition is modeled by (5), (6), and (7).
- *Reward Function:* We seek a policy that maximizes the decision accuracy and minimizes the stopping time $K$. Here, we capture the decision accuracy using the uncertainty associated with each process conditioned on the observations. The uncertainty associated with the $i^{\text{th}}$ process can be quantified using the entropy of its posterior distribution

---

[1]During the training phase, the true value of $\boldsymbol{s}$ is provided, but the optimal selection at each time instant is unknown.

$\begin{bmatrix} \sigma_i(k) & 1 - \sigma_i(k) \end{bmatrix}$. Therefore, the instantaneous reward of the MDP is

$$r(k) = \sum_{i=1}^{N} H(\sigma_i(k-1)) - H(\sigma_i(k)), \quad (10)$$

where $H(x) = -x \log x - (1-x) \log(1-x)$ is the entropy. Then, the long term reward can be defined as the expected discounted reward of the MDP: $\bar{R}(k) = \sum_{l=k}^{K} \gamma^{l-k} r(l)$, where $\gamma \in (0, 1)$ is the discount factor. The discounted reward formulation implies that a reward received $l$ time steps in the future is worth only $\gamma^l$ times what it would be worth if it were received immediately. Thus, this formulation minimizes the stopping time.

Having defined the MDP, we next describe the actor-critic RL algorithm that solves the long-term average reward maximization problem.

### B. Deep Actor-Critic Algorithm

The deep actor-critic algorithm is a deep learning-based RL technique that provides a sequential policy that maximizes the long-term expected discounted reward $\bar{R}(k)$ of a given MDP. The actor-critic framework maximizes the discounted reward using two neural networks: actor and critic networks. The actor learns a stochastic policy that maps the state of the MDP to a probability vector on the set of actions. The critic learns a function that evaluates the policy followed by the actor and gives feedback to the actor. Therefore, the two neural networks interact and adapt to each other.

The components of the actor-critic algorithm are as follows:

*Actor Network:* The actor takes the state of the MDP $\boldsymbol{\sigma}(k-1) \in [0, 1]^N$ as its input. Its output is the probability vector $\boldsymbol{\mu}(\boldsymbol{\sigma}(k-1); \alpha) \in [0, 1]^N$ over the set of processes where $\alpha$ denotes the set of parameters of the actor neural network. The decision maker selects a process $a(k) \sim \boldsymbol{\mu}(\boldsymbol{\sigma}(k-1); \alpha)$, i.e., the $i^{\text{th}}$ process is selected at time $k$ with probability equal to the $i^{\text{th}}$ entry $\mu_i(\boldsymbol{\sigma}(k-1); \alpha)$ of the actor output.

*Reward Computation:* Once the process $a(k)$ is selected, the decision maker receives the corresponding observation $y_{a(k)}$, and the MDP state $\boldsymbol{\sigma}(k-1)$ is updated to $\boldsymbol{\sigma}(k)$ as given by (5). The decision maker also calculates the instantaneous reward $r(k)$ using (10), and the reward value is fed to the critic along with the current and previous states of the MDP.

*Critic Network:* The input to the critic at time $k$ is given by

$$\boldsymbol{\theta}(k) = (\boldsymbol{\sigma}(k), \boldsymbol{\sigma}(k-1), r(k)) \in [0, 1]^N \times [0, 1]^N \times \mathbb{R}.$$

The output of the critic is a scalar critique $\delta(\boldsymbol{\theta}(k); \beta)$ where $\beta$ denotes the set of parameters of the critic neural network. This critique is computed based on the value function $V(\boldsymbol{\sigma}(k))$ of the current MDP state as defined below:

$$V^{\mu}(\boldsymbol{\sigma}) = \mathbb{E}_{a(k) \sim \mu} \left\{ \bar{R}(k) | \boldsymbol{\sigma}(k) = \boldsymbol{\sigma} \right\}.$$

We note that $V^{\mu}(\boldsymbol{\sigma})$ is the expected average future reward when the MDP starts at state $\boldsymbol{\sigma}$ and follows the policy $\mu(\cdot; \theta)$

thereafter. In other words, $V^\mu(\boldsymbol{\sigma})$ indicates the long term desirability of the MDP being in state $\boldsymbol{\sigma}$. The scalar critique takes the form of a temporal difference (TD) error $\delta(\boldsymbol{\theta}(k); \beta)$

$$\delta(\boldsymbol{\theta}(k); \beta) = r(k) + \gamma\hat{V}(\boldsymbol{\sigma}(k)) - \hat{V}(\boldsymbol{\sigma}(k-1)), \quad (11)$$

where $\hat{V}$ is the value function estimate learned by the critic. A positive TD error indicates that the probability of choosing the current action should be increased for the future, and a negative TD error suggests that the probability of choosing $a(k)$ should be decreased.

*Learning Actor Parameters:* The goal of the actor is to choose a policy such that the value function is maximized which in turn maximizes the expected average future reward. Therefore, the actor updates its parameter set $\alpha$ using the gradient descent step by moving in the direction in which the value function is maximized. The update equation for the actor parameters is given by

$$\alpha = \alpha^- + \delta(\boldsymbol{\theta}(k); \beta)\nabla_\alpha[\log\mu_{a(k)}(\boldsymbol{\sigma}(k-1); \alpha)], \quad (12)$$

where $\alpha^-$ is the estimate of the network obtained in the previous time instant [12, Chapter 13].

*Learning Critic Parameters:* The critic chooses its parameters such that it learns the estimate $\hat{V}(\cdot)$ of the state value function $V(\cdot)$ accurately. Therefore, the critic updates its parameter set $\beta$ by minimizing the square of the TD error $\delta^2(\boldsymbol{\theta}(k); \beta)$.

*Termination criterion:* The actor-critic algorithm continues to collect observations until the confidence level on the decision exceeds the desired level $\pi_{\text{upper}}$. We define the confidence level on $\hat{s}_i$ as $\max\{\sigma_i(k), 1 - \sigma_i(k)\}$. Therefore, the stopping criterion is as follows:

$$\min_{i=1,2,\ldots,N} \max\{\sigma_i(k), 1 - \sigma_i(k)\} > \pi_{\text{upper}}. \quad (13)$$

The above components completely describe the actor-critic algorithm, and we next summarize the overall algorithm and discuss its complexity.

## V. OVERALL ALGORITHM

Combing the estimation algorithm in Sec. III and the deep actor-critic method in Sec. IV, we obtain our anomaly detection algorithm. The decision maker collects observations using the selection policy obtained using the actor-critic algorithm until the stopping criterion given in (13) is satisfied. After the actor-critic algorithm terminates, the decision maker computes $\hat{s}$ using (9). We present the pseudo-code of the overall procedure in Algorithm 1.

The computational complexity of our algorithm is determined by the size of the neural networks, the update of the posterior belief vector given by (5)-(7), and the reward computation given by (10). Since all of them have linear complexity in the number of processes $N$, the overall computational complexity of our algorithm is polynomial in $N$. Also, the sizes of all the variables involved in the algorithm are linear in $N$ except for the pairwise conditional probability $\mathbb{P}[s_i|s_j]$

---

**Algorithm 1** Actor-critic RL for anomaly detection

---

**Parameters:** Discount rate $\gamma \in (0,1)$, Upper threshold on confidence $\pi_{\text{upper}} \in (0.5, 1)$
**Initialization:** $\alpha, \beta$ with random weights, $\boldsymbol{\sigma}(0)$ with the prior on each process (can be learned from the training data)
 1: **for** Episode index $= 1, 2, \ldots$ **do**
 2:     Time index $k = 1$
 3:     **repeat**
 4:         Choose a process $a(k) \sim \mu(\boldsymbol{\sigma}(k-1), \alpha)$
 5:         Receive observation $\boldsymbol{y}_{a(k)}(k)$
 6:         Compute $\boldsymbol{\sigma}(k)$ using (5) - (7)
 7:         Compute instantaneous reward $r(k)$ using (10)
 8:         Update the actor neural network using (12)
 9:         Update the critic neural network by minimizing the temporal error $\delta$ in (11) with respect to $\beta$
10:         Increase time index $k = k + 1$
11:     **until** (13) is satisfied
12:     Declare the estimate $\hat{s}$ using (9)
13: **end for**

---

for $i, j = 1, 2, \ldots, N$. Therefore, the memory requirement of the algorithm is $\mathcal{O}(N^2)$. Hence, our algorithm possesses polynomial complexity, unlike the anomaly detection algorithms in [4], [11] that have exponential complexity in $N$. Consequently, our algorithm is more applicable in practical settings.

It is straightforward to extend our algorithm to the case in which the decision maker chooses $n$ processes at a time. In that case, the output layer of the actor has $\binom{N}{n}$ neurons, and we need to update $\boldsymbol{\sigma}(k) \in [0,1]^N$ using the conditional probabilities of the form $\mathbb{P}[s_{i_1}, s_{i_2}, \ldots, s_{i_n}|s_j]$, for $1 < i_1 < i_2 < i_3 < \ldots < N$ and $j = 1, 2, \ldots, N$. Therefore, the overall computational complexity of the resulting algorithm is polynomial in $N$ and the memory requirement is $\mathcal{O}(N^{n+1})$.

## VI. SIMULATION RESULTS

In this section, we empirically study the detection performance of our algorithm. We use two metrics for the performance evaluation: accuracy (the fraction of times the algorithm correctly identifies all the anomalous processes) and stopping time.

### A. Simulation Setup

Our simulation setup is as described below:
*Processes and Their Statistical Dependence:* We consider five processes $N = 5$ and assume that the probability of each process being normal is $q = 0.8$. Here, the first and second processes ($s_1$ and $s_2$) are statistically dependent, and the third and fourth processes ($s_3$ and $s_4$) are also statistically dependent. These pairs of processes are independent of each other and independent of the fifth process ($s_5$). The dependence is captured using the correlation coefficient $\rho \in [0,1]$ that is common to both process pairs:

$$\mathbb{P}[s_1 = s_2 = 0] = \mathbb{P}[s_3 = s_4 = 0] = q^2 + \rho q(1-q)$$
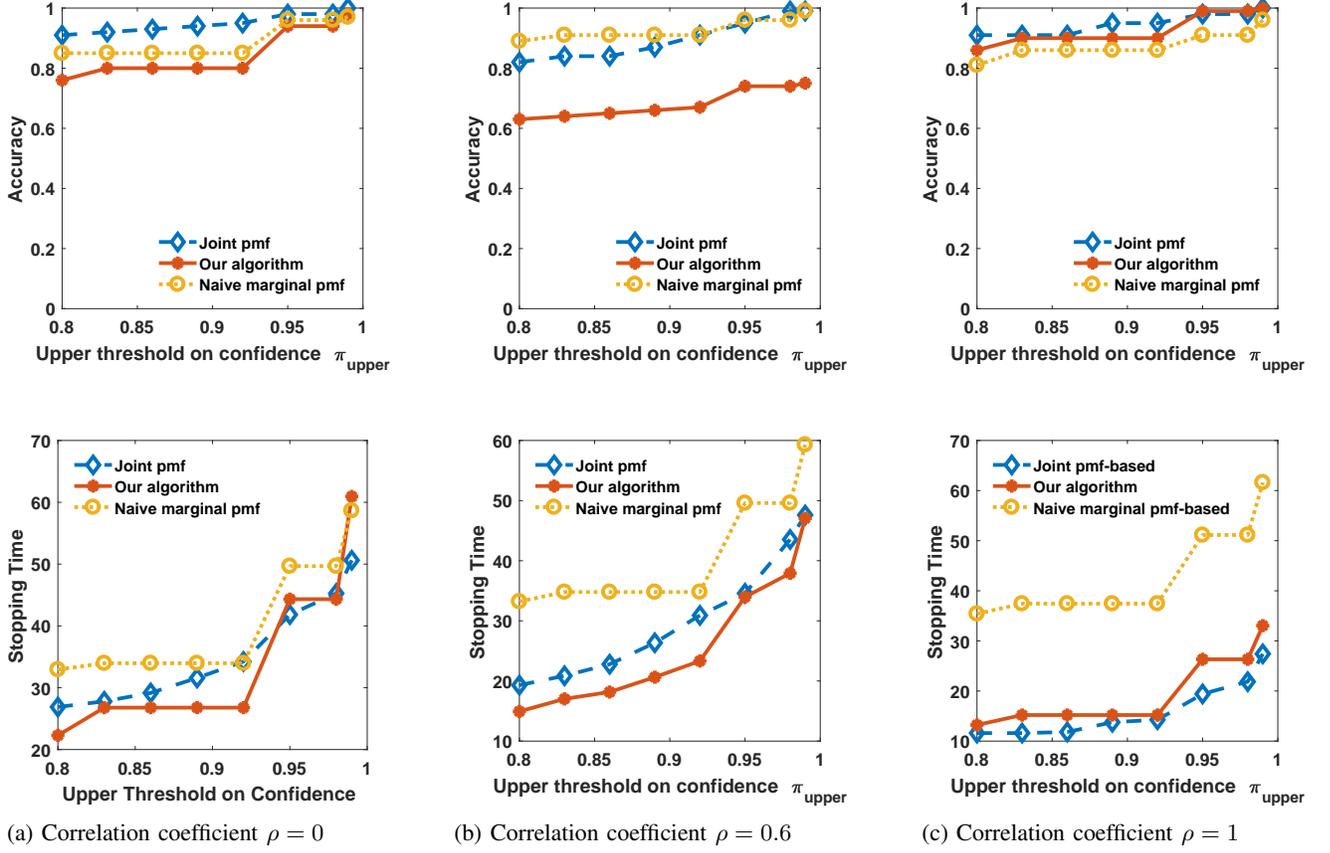$$\mathbb{P}[s_1 \neq s_2] = \mathbb{P}[s_3 \neq s_4] = (1-\rho)q(1-q)$$

Figure 1: Performances of the three different deep actor-critic algorithms as a function of $\pi_{\text{upper}}$.

Also, we assume that the flipping probability $p = 0.2$.

*Implementation of Our Algorithm:* We implement the actor and critic neural networks with three layers and the ReLU activation function between consecutive layers. The output layer of the actor layer is normalized to ensure that $\mu(\cdot)$ is a probability vector over the set of processes. The parameters of the neural networks are updated using the Adam Optimizer, and we set the learning rates of the actor and the critic as $5 \times 10^{-4}$, and $5 \times 10^{-3}$, respectively. Also, we set the discount factor $\gamma = 0.9$.

*Competing Algorithms:* We compare the performance of our algorithm with two other deep actor-critic-based algorithms:

1) *Joint probability mass function (pmf)-based algorithm:* This algorithm refers to the state-of-the-art method for anomaly detection problem presented in [4]. The algorithm is based on the joint posterior probabilities of all the entries of $s \in [0,1]^N$. Since $s$ can take $2^N$ possible values, the complexity of this algorithm is $2^N$. However, the joint probabilities help the algorithm to learn all possible statistical dependencies among the process.

2) *Naive marginal pmf-based algorithm:* We also compare our algorithm with a naive method that also relies on the marginal probabilities $\boldsymbol{\sigma} \in [0,1]^N$. This algorithm is

identical to our algorithm except that at every time instant, this method only updates the entry of $\sigma_{a(k)}(k)$ of $\boldsymbol{\sigma}(k)$ corresponding to the selected process $a(k)$. In other words, this method ignores the possible statistical dependence of the observation $y_{a(k)}(k)$ on the processes other than $a(k)$. Hence, the computational complexity of this algorithm is also $\mathcal{O}(N)$. We note that unlike our algorithm, this algorithm does not use any approximation, and therefore, its updates are always exact.

Our algorithm is a compromise between the above two algorithms and relies on marginal probabilities $\boldsymbol{\sigma}$ while accounting for the possible statistical dependence among the processes.

*B. Discussion of Results*

Our results are summarized in Figs. 1 and 2 and the key inferences from them are as follows:

- The accuracy and the stopping time of all the algorithms increase with $\pi_{\text{upper}}$. This trend is expected due to the fact that as $\pi_{\text{upper}}$ increases, the decision maker requires more observations to satisfy the higher desired confidence level.
- The accuracy of our algorithm is comparable to the other two algorithms when $\rho = 0$ and $\rho = 1$. The accuracy degrades as $\rho$ is close to 0.5. This behavior is because our algorithm uses approximate marginal probabilities to compute the confidence level whereas the other two algorithms use exact
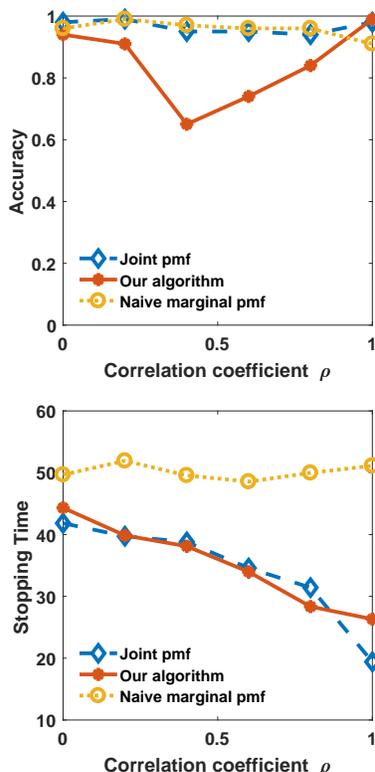
Figure 2: Performances of the three different deep actor-critic algorithms with $\rho$ as a function of $\pi_{\text{upper}} = 0.95$.

values. This approximation in (4) is exact when $\rho = 0$ and $\rho = 1$. As $\rho$ approaches $0.5$, the approximation error increases, and the accuracy decreases.

- The stopping times of the three algorithms are similar when $\rho = 0$. This is because when $\rho = 0$, all the processes are independent. Therefore, the updates of our algorithm are exact. The naive marginal pmf-based algorithm also offers good performance as there is no underlying statistical structure among the processes.

- The stopping times of our algorithm and the joint pmf-based algorithm improve with $\rho$. As $\rho$ increases, the processes become more correlated, and therefore, an observation corresponding to one process has more information about the other correlated processes. However, the naive marginal pmf-based algorithm ignores this correlation and handles the observations corresponding to the different processes independently. Therefore, the stopping time is insensitive to $\rho$. Consequently, the difference between the stopping times of the naive marginal pmf-based algorithm and the other two algorithms increases as $\rho$ increases.

Further, from our experiments, we notice that the average runtime per per process selection decision for the joint pmf-based algorithm, our algorithm, and the naive marginal pmf-based algorithm are 3.2 ms, 2.88 ms, 2.89 ms, respectively. This observation is in agreement with our complexity analysis in Sec. V which implies that the joint pmf-based algorithm is

computationally heavier compared to the other two algorithms. We also recall that the difference between the runtimes of the joint pmf based algorithm and our algorithm grows with $N$.

Thus, we conclude that our algorithm combines the best of two worlds by benefiting from the statistical dependence among the processes (similar to the joint pmf-based algorithm) and offering low-complexity (similar to the naive marginal pmf-based algorithm).

## VII. CONCLUSION

We presented a low-complexity algorithm to detect the anomalous processes among a set of binary processes by observing a single process at a time. The sequential process selection problem was formulated using a Markov decision process whose reward is defined using the entropy of the marginal probabilities of the processes. The optimal process selection policy was obtained via the deep actor-critic algorithm that maximizes the long-term average reward of the MDP. Using numerical results, we established that our algorithm learns and adapts to the underlying statistical dependence among the processes while operating with low complexity. This algorithm relies on approximate marginal probabilities which can lead to performance deterioration when the approximation error is large. A theoretical analysis that quantifies the approximation error is an interesting direction for future work.

## REFERENCES

[1] W.-Y. Chung and S.-J. Oh, "Remote monitoring system with wireless sensors module for room environment," *Sensors Actuators B: Chemical*, vol. 113, no. 1, pp. 64–70, Jan. 2006.

[2] A. Bujnowski, J. Ruminski, A. Palinski, and J. Wtrorek, "Enhanced remote control providing medical functionalities," in *Proc. Inter. Conf. Pervasive Comput. Tech Healthc. Workshops*, May 2013, pp. 290–293.

[3] C. Zhong, M. C. Gursoy, and S. Velipasalar, "Deep actor-critic reinforcement learning for anomaly detection," in *Proc. Globecom*, Dec. 2019.

[4] G. Joseph, M. C. Gursoy, and P. K. Varshney, "Anomaly detection under controlled sensing using actor-critic reinforcement learning," in *Proc. IEEE Inter. Workshop SPAWC*, May 2020.

[5] H. Chernoff, "Sequential design of experiments," *Ann. Math. Stat.*, vol. 30, no. 3, pp. 755–770, Sep. 1959.

[6] S. A. Bessler, "Theory and applications of the sequential design of experiments, k-actions and infinitely many experiments. part i. theory," Stanford Univ CA Applied Mathematics and Statistics Labs, Tech. Rep., 1960.

[7] S. Nitinawarat, G. K. Atia, and V. V. Veeravalli, "Controlled sensing for multihypothesis testing," *IEEE Trans. Autom. Control*, vol. 58, no. 10, pp. 2451–2464, May 2013.

[8] M. Naghshvar and T. Javidi, "Active sequential hypothesis testing," *Ann. Stat.*, vol. 41, no. 6, pp. 2703–2738, 2013.

[9] B. Huang, K. Cohen, and Q. Zhao, "Active anomaly detection in heterogeneous processes," *IEEE Trans. Inf. Theory*, vol. 65, no. 4, pp. 2284–2301, Aug. 2018.

[10] D. Kartik, E. Sabir, U. Mitra, and P. Natarajan, "Policy design for active sequential hypothesis testing using deep learning," in *Proc. Allerton*, Oct. 2018, pp. 741–748.

[11] G. Joseph, C. Zhong, M. C. Gursoy, S. Velipasalar, and P. K. Varshney, "Anomaly detection under controlled sensing using actor-critic reinforcement learning," in *Proc. IEEE Globecom*, Dec. 2020.

[12] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.