# Enhancement of a state-of-the-art RL-based detection algorithm for Massive MIMO radars

**FRANCESCO LISI**, Student Member, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy

**STEFANO FORTUNATI**, Senior Member, IEEE
Université Paris-Saclay, CNRS, CentraleSupélec, L2S, Gif-sur-Yvette & DR2I-IPSA, Ivry-sur-Seine, France

**MARIA SABRINA GRECO**, Fellow, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy

**FULVIO GINI**, Fellow, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy

*Abstract—* **In the present work, a reinforcement learning (RL) based *adaptive* algorithm to optimise the transmit beampattern for a co-located massive MIMO radar is presented. Under the massive MIMO regime, a robust Wald-type detector, able to guarantee certain detection performances under a wide range of practical disturbance models, has been recently proposed. Furthermore, an RL/cognitive methodology has been exploited to improve the detection performance by learning and interacting with the surrounding unknown environment. Building upon previous findings, we develop here a fully adaptive and *data-driven* scheme for the selection of the hyper-parameters involved in the RL algorithm. Such an adaptive selection makes the Wald-RL-based detector independent of any ad-hoc, and potentially sub-optimal, manual tuning of the hyper-parameters. Simulation results show the effectiveness of the proposed scheme in harsh scenarios with strong clutter and low SNR values.**

*Index Terms—* **Adaptive selection, beamforming, constant false alarm rate, massive MIMO radar, Reinforcement Learning, SARSA, target detection.**

F. Lisi, M. S. Greco and F. Gini are with Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy (e-mail: francesco.lisi@phd.unipi.it, maria.greco@unipi.it, fulvio.gini@unipi.it). S. Fortunati is with Université Paris-Saclay, CNRS, CentraleSupélec, Laboratoire des signaux et systèmes, 91190, Gif-sur-Yvette & DR2I-IPSA, 94200, Ivry-sur-Seine, France (e-mail: stefano.fortunati@centralesupelec.fr).

The code related to this work can be found at https://github.com/lisifra96/Improved_RL_algorithm_mMIMO_radar

Color versions of one or more of the figures in this article are available online at http://ieeexplore.ieee.org.

## I. INTRODUCTION

The concept of Cognitive Radar (CR) has been firstly introduced by Haykin [1]. The need of adaptation to the changes in the environment is intrinsic in the radar detection problem due to the presence of multiple sources of non-stationarity, such as variations in the clutter statistics over time or changes in the target scenario. Unlike communication systems where the transmitter and receiver are physically separated, in a monostatic radar system the transmitter and the receiver are in the same position, allowing the latter to easily transmit information to the former. MIMO radar systems can be divided into two fundamentals categories: widely separated radars and co-located radars. As [2] suggests, a radar with widely separated antennas can exploit the spatial diversity of the target radar cross section to obtain a diversity gain similar to the one in communications. In a co-located MIMO radar, the antennas are closely spaced and each element of the array transmits a different probing signal, contrary to standard phased array where all the elements transmit the same waveform with variable amplitude and phase. In [3] the authors discuss the advantages of parameter identifiability led by co-located MIMO radars. Following the recent works [4], [5], in this correspondence, we focus our attention on co-located massive MIMO radars exploiting RL-based techniques [6], [7] to implement the *cognition loop* [8]. Specifically, in [4] the authors proved that, if the number of virtual spatial antenna channels is high enough (massive MIMO regime), a robust Wald-type test can be derived to guarantee the constant false alarm rate (CFAR) property under a wide variety of disturbance models using a single snapshot. Then, in [5], a reinforcement learning (RL) approach is proposed to optimise the transmission beampattern of a co-located massive MIMO radar exploiting the robust Wald-type detector in [4].

Two original contributions are proposed in this correspondence. The first one is the introduction of two new policies, called *quasi $\varepsilon$-greedy policy* and *quasi $\varepsilon$-greedy policy with target recovery*. The former improves the performance of the system by forcing the RL algorithm to focus its power in a number of angular bins greater or equal to the number of detected targets; in addition, the latter exploits a mechanism that allows the system to recover a missed target more quickly. The second original contribution consists of an adaptive algorithm that selects the SARSA algorithm hyper-parameters ($\varepsilon$ and $\alpha$) based on the received signal. This allows the system to adapt to the variations of the surrounding environment and to achieve better detection capability both in stationary and non-stationary scenarios. Afterwards, the RL-based algorithm proposed in [5] with the quasi $\varepsilon$-greedy policy with target recovery and the adaptive selection of $\varepsilon$ and $\alpha$ is tested in various simulation scenarios. The results confirmed that the updated algorithm achieves higher detection performance, compared to the one in [5]. Both

our algorithm and the one proposed in [5] exploit the robust Wald-type detector described in [4].

In Section II we describe the detection problem and the main properties of the robust Wald-type detector derived in [4]. Section III provides the background on the RL-based detector proposed in [5]. In Section IV two variations of the $\varepsilon$-greedy policy are proposed, while an adaptive algorithm to select the $\varepsilon$ and $\alpha$ parameters is described in Section V. In Section VI we validate the proposed contributions via simulation. The simulated scenarios are described in Appendix A. Finally, conclusions are drawn in Section VII.

**NOTATIONS:** In this paper we use upper case letters $\boldsymbol{A}$ and lower case ones $\boldsymbol{a}$ to denote matrices and vectors respectively. $(\cdot)^T$ and $(\cdot)^H$ denote a matrix transpose and conjugate transpose respectively, while $(\cdot)^*$ denotes the conjugate operator. $\boldsymbol{I}_N$ denotes the $N \times N$ identity matrix, while $\boldsymbol{0}_N$ denotes an all zeros $N \times N$ matrix. $E\{\cdot\}$ denotes the statistical expectation. The Kronecker product is represented by $\otimes$. A closed interval of numbers between $a$ and $b$ is denoted by $[a, b]$, while a set containing only $a$ and $b$ is denoted by $\{a, b\}$. The absolute value is represented by $|\cdot|$. The relative complement of set $\mathcal{A}$ with respect to set $\mathcal{B}$ is denoted as $\mathcal{B}\backslash\mathcal{A}$.

## II. THE CFAR DETECTION ALGORITHM

Consider a co-located MIMO radar with $N_T$ transmit and $N_R$ receive antennas. Both transmit and receive arrays are uniform linear arrays (ULA) with inter-element distance $d = \lambda/2$. The transmitted signal vector can be expressed as $\boldsymbol{s}(t) = \boldsymbol{W}\boldsymbol{\Phi}(t)$ where $\boldsymbol{W} \in \mathbb{C}^{N_T \times N_T}$ is a weighting matrix and $\boldsymbol{\Phi}(t) \in \mathbb{C}^{N_T}$ is a vector containing a set of orthonormal signals. After sampling the correct range-doppler bin at the output of the matched filter the received signal can be expressed as [4]

$$\boldsymbol{y} = \alpha\boldsymbol{h} + \boldsymbol{c} \in \mathbb{C}^{N \times 1}, \tag{1}$$

with $N = N_T N_R$ and $\boldsymbol{h} = \boldsymbol{W}^T\boldsymbol{a}_T(\nu_0) \otimes \boldsymbol{a}_R(\nu_0)$. Finally, $\boldsymbol{a}_T(\nu_0)$ and $\boldsymbol{a}_R(\nu_0)$ are the transmit and receive steering vectors that depend on the spatial frequency $\nu_0 = (d/\lambda)\sin(\theta_0)$ where $\theta_0$ is the target angle of arrival (in a given reference frame).

Assuming that the received signal is processed at each time instant $k$ by a bank of $L$ spatial filters tuned to a fixed grid of angular bins $\Theta = \{\theta_l\}_{l=1}^L$, the detection problem can be formulated as the following hypothesis testing problem [5]

$$\begin{aligned} \mathcal{H}_0 &: \boldsymbol{y}_{k,l} = \boldsymbol{c}_{k,l}, \\ \mathcal{H}_1 &: \boldsymbol{y}_{k,l} = \alpha_{k,l}\boldsymbol{h}_{k,l} + \boldsymbol{c}_{k,l}. \end{aligned} \tag{2}$$

In the following, we assume that the grid is chosen in order to uniformly span the spatial frequency interval $[-0.5, 0.5]$.

As a test statistic, the robust Wald-type detector $\Lambda_{k,l}$ is adopted [4], [9]:

$$\Lambda_{k,l} = \frac{2|\boldsymbol{h}_{k,l}^H\boldsymbol{y}_{k,l}|^2}{\boldsymbol{h}_{k,l}^H\hat{\boldsymbol{\Gamma}}_{k,l}\boldsymbol{h}_{k,l}}, \tag{3}$$

where $\hat{\boldsymbol{\Gamma}}_{k,l}$ is an estimate of the covariance matrix of the noise vector $\boldsymbol{c}_{k,l}$. The detector discriminates between $\mathcal{H}_0$ and $\mathcal{H}_1$ by comparing the statistic with a threshold $\lambda$ in each angular bin. When $N \to \infty$, the statistic satisfies the CFAR property, i.e. the Probability of False Alarm ($P_{FA}$) is constant and can be selected by choosing the threshold as $\lambda = -2\ln(P_{FA})$, under a wide rage of practical disturbance processes with unknown statistical characterisation. For a $P_{FA}$ equal to $10^{-4}$, if $N \geq 10^4$ the MIMO radar is assumed to operate in the massive MIMO regime and the previous property is satisfied [4].

## III. RL-BASED ALGORITHM

After having ensured the CFARness of the detection scheme through the Wald-type detector in (3), the work [5] focused on how to exploit the large degrees of freedom, offered by a massive MIMO radar, to maximise the Probability of Detection ($P_D$). This can be achieved by properly selecting the weighting matrix $\boldsymbol{W}_k$ to shape the transmit beampattern $\boldsymbol{a}_T^T(\nu)\boldsymbol{W}_k\boldsymbol{W}_k^H\boldsymbol{a}_T^*(\nu)$ [10]. In order to do so, an RL-based algorithm has been proposed to allow the radar focusing the power in the angular bins associated to the targets' angular position [5]. In the following, a concise summary of the main definitions and results obtained in [5] will be provided since they represent the starting point of the original developments proposed in this correspondence. Let us start by introducing the Markov Decision Process (MDP) characterising the learning/detection task at hand. For the sake of clarity, it is worth mentioning that the following definitions of state, action and reward, associated to the considered MDP, are slightly different with respect to the ones introduced in [5]. Even though the essence remains exactly the same, this new formulation is more precise and self-consistent.

### A. The set of the states

The state space of our MDP is denoted as $\mathcal{S} \triangleq \{s^{(i)}\}_{i=0}^K$, where $K < \infty$ is the maximum number of detectable target.

The state of the system at time instant $k$ is defined as $s_k = s^{(i_k)}$ with [5]

$$i_k \triangleq \min\left\{\sum_{l=0}^{L-1}\bar{\Lambda}_{k,l}, K\right\}, \tag{4}$$

where $\bar{\Lambda}_{k,l} \triangleq u(\Lambda_{k,l} - \lambda)$ and $u(\cdot)$ corresponds to the Heaviside step function [1].

### B. The set of the actions

The action set of the MDP is indicated as $\mathcal{A} \triangleq \{a^{(j)}\}_{j=0}^K$. If the action $a^{(j_k)}$ is selected at time

---

[1] In this correspondence the notation $(\cdot)_{k,l}$ is used, unlike the one $(\cdot)_l^k$ used in [5] .

instant $k$, then the beamforming algorithm focuses its power in the $j_k$ angular bins with the highest decision statistic.

Let $\{l_k^{(n)}\}_{n=1}^{L}$ be the set containing the angular bin's indexes sorted in descending order with respect to the decision statistic at time instant $k$ ($\Lambda_{k,l_k^{(1)}} \geq ... \geq \Lambda_{k,l_k^{(L)}}$), we define the set

$$\Omega_k \triangleq \begin{cases} \emptyset\,, & j_k = 0 \\ \{l_k^{(n)}\}_{n=1}^{j_k}\,, & j_k \neq 0 \end{cases}, \quad (5)$$

containing the indexes of the angular bins associated to the $j_k$ highest decision statistics.

The system chooses the matrix $\boldsymbol{W}_k$ according to [5], [11]

$$\boldsymbol{W}_k = \begin{cases} \boldsymbol{W}_{ort} \triangleq \sqrt{\frac{P_{max}}{N_T}} \cdot \boldsymbol{I}_{N_T}\,, & \Omega_k = \emptyset \\ \underset{\boldsymbol{W}}{\arg\max}\,\underset{l \in \Omega_k}{\min} \quad \boldsymbol{a}_T^T(\nu_l)\boldsymbol{W}\boldsymbol{W}^H\boldsymbol{a}_T^*(\nu_l) \\ \text{subject to} \qquad \text{tr}\{\boldsymbol{W}\boldsymbol{W}^H\} \leq P_{max} \end{cases}, \quad \Omega_k \neq \emptyset \quad (6)$$

where $P_{max}$ is the maximum transmitted power.

## C. The reward function

Let's define the set $\Psi_k \triangleq \{l_k^{(n)}\}_{n=1}^{K}$, where $\{l_k^{(n)}\}_{n=1}^{L}$ is the same set as in Section B, and the sets

$$\Phi_k \triangleq \begin{cases} \emptyset\,, & i_k = 0 \\ \{l_k^{(n)}\}_{n=1}^{i_k}\,, & i_k \neq 0 \end{cases}, \quad (7)$$

and $\bar{\Phi}_k \triangleq \Psi_k \backslash \Phi_k$, where $i_k$ is defined in (4).

The reward is defined as [5]

$$r_k \triangleq \sum_{l \in \Phi_k} \hat{P}_{D,k,l} - \sum_{l \in \bar{\Phi}_k} \hat{P}_{D,k,l} \quad (8)$$

where $\hat{P}_{D,k,l}$ is an estimate of the $P_D$ of a target located in the $l^{th}$ angular bin at the $k^{th}$ time instant.

## D. SARSA algorithm

The goal of RL algorithms is to find the optimal policy $\pi^*$, i.e. the one that maximises the *state value function* $V_\pi(s) \triangleq E_\pi \left\{ \sum_{h=0}^{+\infty} \gamma^h r_{k+h+1} \big| s_k = s \right\}$, $\forall s \in \mathcal{S}$, where $\gamma \in [0,1]$ is a damping factor. By defining the *state-action value function* associated to the policy $\pi$ as $Q(s,a) \triangleq E_\pi \left\{ \sum_{h=0}^{+\infty} \gamma^h r_{k+h+1} \big| (s_k = s) \cap (a_k = a) \right\}$, it can be proved that the greedy action associated to the state $s$ is equal to $\pi^*(s) = \arg\max_{a' \in \mathcal{A}} Q_{\pi^*}(s, a')$ [6], [7]. The SARSA algorithm, named after the update rule sequence *State-Action-Reward-State-Action*, allows the system to recursively compute the $\boldsymbol{Q}$ matrix associated to the optimal policy when the dynamics of the environment are unknown.

The algorithm proposed in [5] consists of setting the initial values $s_0 = s^{(0)}$, $a_0 = a^{(0)}$ and $\boldsymbol{Q}_0 = \boldsymbol{0}_{K+1}$ and proceeding by computing the new state $s_{k+1}$ and reward $r_{k+1}$, selecting a new action $a_{k+1}$ following the $\varepsilon$-greedy policy and then updating $\boldsymbol{Q}$ as

$$Q_{k+1}(s_k, a_k) = Q_k(s_k, a_k) + $$
$$+ \alpha_{k+1}(r_{k+1} + \gamma Q_k(s_{k+1}, a_{k+1}) - Q_k(s_k, a_k)). \quad (9)$$

Once the action has been selected the algorithm computes the $\boldsymbol{W}_{k+1}$ matrix following (6) and then the system transmits the new pulses. The *learning rate* $\alpha_k \in (0,1)$ is the weight given to the new information with respect to the old one.

## E. $\varepsilon$-greedy policy

To ensure the convergence of the SARSA algorithm, the new action $a_{k+1}$ must be selected according to a policy that guarantees that each state-action pair is visited infinitely many times. The $\varepsilon$-greedy policy is one of the most used ones in the RL-literature [6], [7]. The new action is selected as

$$\pi_k^{(1)}(s_k) = \begin{cases} a_k^{(greedy)}, & \text{w.p. } 1 - \varepsilon \\ \mathbb{U}\left\{\mathcal{A} - \{a_k^{(greedy)}\}\right\}, & \text{w.p. } \varepsilon \end{cases}, \quad (10)$$

where "w.p." stands for *with probability*, $a_k^{(greedy)} = \arg\max_{a' \in \mathcal{A}} Q_k(s_k, a')$ and $\mathbb{U}\{\mathcal{C}\}$ denotes a function that selects randomly one of the elements in the set $\mathcal{C}$ with uniform probability. The $\varepsilon$ parameter controls the *exploration-exploitation tradeoff*: low $\varepsilon$ values correspond to a system that chooses the greedy action most of the time (exploitation), while high values to one that selects a random action more frequently (exploration).

## IV. Policy improvement

This section and the following Section V present the original contribution of this correspondence. We start by providing two improved versions of the $\varepsilon$-greedy policy.

## A. Quasi $\varepsilon$-greedy policy

The problem arising from choosing the $\varepsilon$-greedy policy in our specific application can be easily explained with an example.

*Example 1.* Consider a static scenario with two targets as in Scenario 1 (see Appendix A). Suppose that the system performs the optimal action at the $k^{th}$ iteration, i.e. selects action $a_k = a^{(2)}$ and the set $\Omega_k$, defined in (5), contains the indices associated to the actual position of the two targets. If the system chooses the random action $a_{k+1} = a^{(1)}$ at the next time instant and focuses all the power in only one of the two targets, with high probability it will end up in $s_{k+2} = s^{(1)}$. At this point, even if the system selects the action $a_{k+2} = a^{(2)}$, the set $\Omega_{k+2}$ isn't guaranteed to contain the indexes of the two targets. The system will stay in state $s^{(1)}$ until the missed target will be in the set $\Omega_{k+m}$.

To overcome this issue a new policy is proposed here. We define *the quasi $\varepsilon$-greedy policy* as:

$$\pi_k^{(2)}(s_k) = \begin{cases} \mathbb{U}\left\{\mathcal{A}'(s_k) - \{a_k^{(greedy)}\}\right\}, & \text{w.p. } \varepsilon \\ a_k^{(greedy)}\,, & \text{w.p. } 1 - \varepsilon \end{cases} \quad (11)$$

where $\mathcal{A}'(s^{(i)}) \triangleq \{a^{(j)}, j = i, ..., K\}$. If the quasi $\varepsilon$-greedy policy is combined with $\boldsymbol{Q}_0 = \boldsymbol{I}_{K+1}$ then, when

the system is in state $s^{(i)}$, it can't focus its power in a number of angular bins less than $i$, that corresponds to the number of detected targets if these aren't more than $K$.

### B. Quasi $\varepsilon$-greedy policy with target recovery

Although the quasi $\varepsilon$-greedy policy solved one of the main issues encountered with the $\varepsilon$-greedy policy, it doesn't solve the *target loss problem* described in the following example.

*Example 2*. Consider a scenario with two targets in the radar scene, like in Scenario 1 (see Appendix A). Suppose that, at a given time instant $k$, the system is in state $s^{(2)}$, having detected both targets, and randomly selects action $a^{(5)}$. Although two of the five angular bins in which the system focuses its power correspond to the position of the targets, the system might lose one of them if its SNR is low and end up in state $s^{(1)}$. Now, with high probability, the two highest decision statistics still correspond to the angular position of the targets. Consequently, if the system selects action $a^{(2)}$, it recovers the lost target and gets back in state $s^{(2)}$. On the other hand, since the system is in state $s^{(1)}$, it might happen that it selects action $a^{(1)}$ randomly or because this is the greedy action associated to state $s^{(1)}$, especially in the early phases when the system hasn't figured out the scene yet. In this case the system focuses all its power in the direction of the detected target and the decision statistic of the missed one drops. Now it will take a long time for the system to recover it.

To solve this "missed target" problem, we define a new policy as:

$$\pi_k^{(3)}(s_k, s_{k-1}) \triangleq \begin{cases} \arg\max_{a \in \mathcal{A}} Q_k(s_{k-1}, a) \,, & i_k < i_{k-1} \\ \pi_k^{(2)}(s_k) \,, & i_k \geq i_{k-1} \end{cases} . \tag{12}$$

In words, this policy can be motivated as follows: when the system detects an higher or equal number of targets, it chooses the new action according to the quasi $\varepsilon$-greedy policy. On the contrary, if the number of detected targets is smaller than the one at the previous iteration the system tries to recover them as soon as possible by choosing the greedy action associated to the state at the previous time instant $s_{k-1}$.

### V. Adaptive selection of $\varepsilon$ and $\alpha$

The $\varepsilon$ and $\alpha$ hyper-parameters are both kept constant in the algorithm proposed in [5]. Although this is one of the most popular choices for non-stationary environments, it has two major drawbacks. Firstly, in order to use the algorithm, the user must select them in advance and this operation requires a certain knowledge of both the problem and algorithm, thus limiting its usability. Secondly, the parameters' optimal value changes over time, both in stationary and non-stationary environments. To overcome these problems we propose an adaptive algorithm that selects both the parameters based on the collected data.

TABLE I: Adaptive $\varepsilon$ and $\alpha$ algorithm parameters. $\varepsilon_{min}$ ($\alpha_{min}$) and $\varepsilon_{max}$ ($\alpha_{max}$) are the endpoints of the interval to which $\varepsilon$ ($\alpha$) belongs; $c_1$ and $c_2$ correspond to the multiplicative decrease and increase factors; $\eta_1$ and $\eta_2$ are the lower and upper threshold respectively.

| $x$ | $x_{min}$ | $x_{max}$ | $c_1$ | $c_2$ | $\eta_1$ | $\eta_2$ |
|---|---|---|---|---|---|---|
| $\varepsilon$ | 0.1 | 0.8 | 0.8 | 2 | 0.5 | 1.8 |
| $\alpha$ | 0.2 | 0.6 | 0.9 | 2.5 | 0.5 | 1.8 |

Let $r_k$ be the reward at time instant $k$ and $d_k$ the sequence defined as

$$d_k \triangleq \begin{cases} r_1 & , \ k = 1 \\ r_k - r_{k-1} & , \ k \neq 1 \end{cases} . \tag{13}$$

Then, $\varepsilon$ and $\alpha$ can be updated, at the $k^{th}$ iteration, according to the following strategy:

$$x_{k+1} = \begin{cases} \max\{c_1 \cdot x_k \,, \ x_{min}\}, & |d_k| < \eta_1 \\ \min\{c_2 \cdot x_k \,, \ x_{max}\}, & \eta_1 < \ |d_k| < \eta_2 \\ x_{max}, & |d_k| > \eta_2 \end{cases} \tag{14}$$

where $x$ corresponds to $\varepsilon$ or $\alpha$, $c_1 \in (0,1)$ and $c_2 \in (1, +\infty)$. The initial value of $x$ is set to $x_0 = x_{max}$. $c_1$, $c_2$, $\eta_1$ and $\eta_2$ are constants, but differ for the $\varepsilon$ and $\alpha$ algorithm. Table I lists all the constants' values. $x_{k+1}$ is not updated if the system was in exploration mode in the two previous time instants, i.e. $k-1$ and $k$. If the system selects a random action at time instant $k$, the reward $r_{k+1}$ may drop due to a bad choice of the action causing $|d_{k+1}| = |r_{k+1} - r_k|$ to surpass $\eta_1$ even though the scenario hasn't changed. If the system then chooses the correct action at time $k+1$, the reward $r_{k+2}$ rises back to a value around $r_k$, but $|d_{k+2}| = |r_{k+2} - r_{k+1}|$ is likely to be over $\eta_1$ due to the low value of $r_{k+1}$.

The thresholds $\eta_1$ and $\eta_2$ depend on the definition of the reward. For the one defined in (8), if the system misses one target or detects a new one the reward variation will be around 1 depending on the target SNR. A good choice for $\eta_1$ is 0.5. Moreover, $\eta_2$ should be chosen high enough to guarantee that the value is reset only when a sudden change in the scenario happens: some empirical tests suggested that $\eta_2 \geq 1.8$ meets this specification.

### VI. SIMULATION RESULTS

In this section the performances of the two policies and the adaptive algorithm previously described are validated via simulations in Scenarios 1 and 2 (see Appendix A). Then the RL-based cognitive algorithm (RL-C) with the adaptive selection of $\varepsilon$ and $\alpha$ and the quasi $\varepsilon$-greedy policy with target recovery is compared with the orthogonal algorithm and a "non RL-based" cognitive algorithm (NRL-C) in two non-stationary scenarios (3 and 4). The orthogonal algorithm is the one with an omnidirectional beampattern, while the NRL-C algorithm focuses the power in the angular bins where a detection occurred. As performance benchmark, we plot an upper bound on

TABLE II: Simulation parameters. $L$ is the number of spatial filters in the receiver and $K$ is the maximum number of detectable targets; $\gamma$ is the damping factor of the cumulative reward and $\boldsymbol{Q}_0$ is the initial value of the *state-action value* matrix.

| $P_{FA}$ | $N_T$ | $N_R$ | L | K | $P_{max}$ | $\gamma$ | $\boldsymbol{Q}_0$ |
|---|---|---|---|---|---|---|---|
| $10^{-4}$ | 100 | 100 | 20 | 5 | 1 | 0.8 | $\boldsymbol{I}_{K+1}$ |

TABLE III: $P_D$ of the four targets in Scenario 3 in [5] for a fixed $P_{FA}$ value of $10^{-4}$. The values in the second row are obtained using the proposed RL-C algorithm, while the ones in the third row are extracted from Figure 6 in [5].

| Target | 1 ($\nu = -0.2$) | 2 ($\nu = 0$) | 3 ($\nu = 0.2$) | 4 ($\nu = 0.3$) |
|---|---|---|---|---|
| RL-C | 1.00 | 0.98 | 0.99 | 0.97 |
| [5] | 1.00 | 0.74 | 0.91 | 0.73 |

the $P_D$ obtained by a clairvoyant beamformer that focuses its power in the exact (and generally unknown) direction of the targets. Finally, the proposed RL-C algorithm is compared with the one in [5] in a stationary scenario. Due to a lack of space only the most relevant figures are shown, but the interested reader can obtain all the simulations' results using the code available via the link in the first page of the correspondence. Table II lists all the parameters' values used in the simulations.

Figure 1 shows the $P_D$ of target 2 in Scenario 1 for the $\varepsilon$-greedy policy and the two proposed policies. The figure shows that both the quasi $\varepsilon$-greedy policy and quasi $\varepsilon$-greedy policy with target recovery have better detection performances than the classical $\varepsilon$-greedy one. In particular, the quasi $\varepsilon$-greedy policy with target recovery has the best performances among the three policies.

Figure 2 compares the performance of the adaptive $\varepsilon$ algorithm with the two static cases with $\varepsilon = \varepsilon_{min}$ and $\varepsilon = \varepsilon_{max}$ while keeping $\alpha$ constant ($\alpha = 0.5$); Figure 3 does the same for $\alpha$ when $\varepsilon$ is kept constant ($\varepsilon = 0.5$). All the graphs show that the adaptive algorithm combines the positive effects of having high values of both the parameters in the initial transitory phase, when the system has to gather information about the surrounding environment, and low values in the following phase when the system approaches a steady state.

Scenario 3 consists of two targets with fixed angular position and variable SNR, which is shown in the inset in Figure 4a. Figure 4 shows the $P_D$ of both targets. Since the SNR increases in the interval [1,100] and decreases in [101,200], the $P_D$ varies accordingly. It is interesting to note that there is a delay between the drop of the $P_D$ and the time instant when the SNR starts to go down ($k = 100$) for both the RL and NRL cognitive algorithms. For example, considering the performance of the RL-C algorithm in Figure 4b, the $P_D$ goes down after $k = 110$; in the interval [101,110] the positive effect of the algorithm's learning capability prevails over the negative effect of the decreasing of the target's SNR. The RL-based algorithm shows better detection performance than the non RL-based cognitive one, confirming that the system is able to exploit the information gathered from the surrounding environment. Since the SNR of both targets is low, the orthogonal algorithm can't detect them. Even though the RL-based algorithm shows far better performance than the other algorithms, there's still a gap with the upper bound on the $P_D$.

Scenario 4 consists of three stationary targets. Target 1 and 2 are in the scene at time $k = 1$ and disappear at time $k = 101$ and $k = 301$ respectively. Target 3 appears at time $k = 201$ and stays in the scene until the end of the simulation. Figure 5 confirms that the RL-based algorithm is the one with the best detection performance among the three algorithms. Figure 6 and Figure 7 show the temporal evolution of the $\varepsilon$ and $\alpha$ parameters: the adaptive algorithm is able to track the variations in the scene and adjust the parameters accordingly.

Finally, Table III compares the performance of the proposed RL-C algorithm with one obtained in [5]. Thanks to the developed quasi $\varepsilon$-greedy policy with target recovery and the adaptive selection of $\varepsilon$ and $\alpha$, a significant improvement has been achieved.

## VII. CONCLUSION

In the present paper we introduced two variations of the $\varepsilon$-greedy policy. In addition, we proposed an adaptive algorithm to select the SARSA $\varepsilon$ and $\alpha$ parameters, which increases both the RL-based algorithm performance and usability, since the user doesn't have to set them manually. Then the enhanced version of the algorithm proposed in [5] was tested in various stationary and non-stationary scenarios. The results confirmed that the updated algorithm leads to a significant improvement of the detection capability of the massive MIMO radar system. On the other hand, the new policies and the adaptive selection of the hyper-parameters slightly increase the computational overhead, which must be taken into account when the algorithm is implemented by a real-time system. Furthermore, the time complexity of the algorithm grows exponentially with $N_T$, due to the growth of the dimension of the weighting matrix $W$ [5], limiting its scalability. Future works will investigate the possibility to fill the gap between the performance of the improved RL-based algorithm and the upper bound. More specifically, we will try to fuse the position information (i.e. target tracking) with the RL-based detection algorithm.

## Appendix A
## SIMULATION FRAMEWORK DESCRIPTION

In all the simulated scenarios the noise process is an AR(6) process with t-distributed independent and

identically distributed (i.i.d.) innovations with the same parameters as the ones in [5].

Table IV contains all the details that characterise the four simulated scenarios.

TABLE IV: Target scenarios.

| Scenario | Time Interval | Target | Angular Bin | $\nu$ | $SNR_{dB}$ |
|---|---|---|---|---|---|
| 1 | [1,300] | 1 | 7 | -0.20 | -20 |
|   |         | 2 | 16 | 0.25 | -20 |
| 2 | [1,100] | 1 | 17 | 0.30 | -20 |
| 3 | [1,200] | 1 | 7 | -0.20 | variable |
|   |         | 2 | 16 | 0.25 | (Inset Fig.4a) |
| 4 | [1,100] | 1 | 5 | -0.30 | -18 |
|   |         | 2 | 13 | 0.10 | -21 |
|   | [101,200] | 2 | 13 | 0.10 | -21 |
|   | [201,300] | 2 | 13 | 0.10 | -21 |
|   |           | 3 | 17 | 0.30 | -20 |
|   | [301,400] | 3 | 17 | 0.30 | -20 |

REFERENCES

[1] S. Haykin, "Cognitive radar: a way of the future," *IEEE signal processing magazine*, vol. 23, no. 1, pp. 30–40, 2006.

[2] A. M. Haimovich, R. S. Blum, and L. J. Cimini, "MIMO radar with widely separated antennas," *IEEE Signal Processing Magazine*, vol. 25, no. 1, pp. 116–129, 2007.

[3] J. Li and P. Stoica, "MIMO radar with colocated antennas," *IEEE Signal Processing Magazine*, vol. 24, no. 5, pp. 106–114, 2007.

[4] S. Fortunati, L. Sanguinetti, F. Gini, M. S. Greco, and B. Himed, "Massive MIMO radar for target detection," *IEEE Transactions on Signal Processing*, vol. 68, pp. 859–871, 2020.

[5] A. M. Ahmed, A. A. Ahmad, S. Fortunati, A. Sezgin, M. Greco, and F. Gini, "A reinforcement learning based approach for multi-target detection in Massive MIMO radar," *IEEE Transactions on Aerospace and Electronic Systems*, pp. 1–1, 2021.

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

[7] M. Mohri, A. Rostamizadeh, and A. Talwalkar, *Foundations of machine learning*, ch. 17. MIT press, 2018.

[8] M. S. Greco, F. Gini, P. Stinco, and K. Bell, "Cognitive radars: On the road to reality: Progress thus far and possibilities for the future," *IEEE Signal Processing Magazine*, vol. 35, no. 4, pp. 112–125, 2018.

[9] S. Fortunati, L. Sanguinetti, F. Gini, M. S. Greco, and B. Himed, "Erratum to "Massive MIMO radar for target detection"," *IEEE Transactions on Signal Processing*, vol. 69, pp. 3235–3235, 2021.

[10] D. R. Fuhrmann and G. San Antonio, "Transmit beamforming for MIMO radar systems using signal cross-correlation," *IEEE Transactions on Aerospace and Electronic Systems*, vol. 44, no. 1, pp. 171–186, 2008.

[11] L. Wang, S. Fortunati, M. S. Greco, and F. Gini, "Reinforcement learning-based waveform optimization for MIMO multi-target detection," in *2018 52nd Asilomar Conference on Signals, Systems, and Computers*, pp. 1329–1333, IEEE, 2018.

Université Paris-Saclay, CNRS, CentraleSupélec, L2S, Gif-sur-Yvette & DR2I-IPSA, Ivry-sur-Seine, France

MARIA SABRINA GRECO, Fellow, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy

FULVIO GINI, Fellow, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy

FRANCESCO LISI, Student Member, IEEE
Dipartimento di Ingegneria dell'Informazione, Università di Pisa, Italy
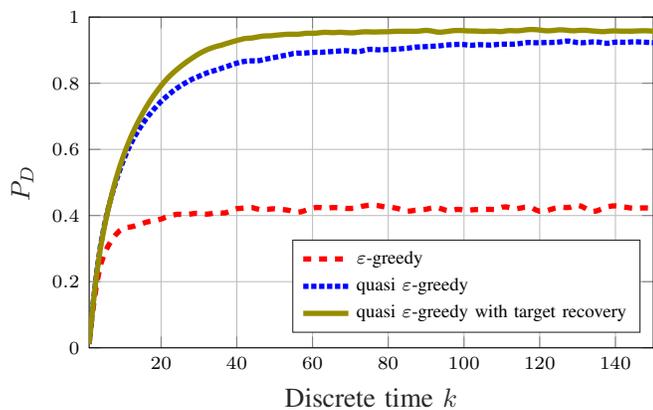
STEFANO FORTUNATI, Senior Member, IEEE

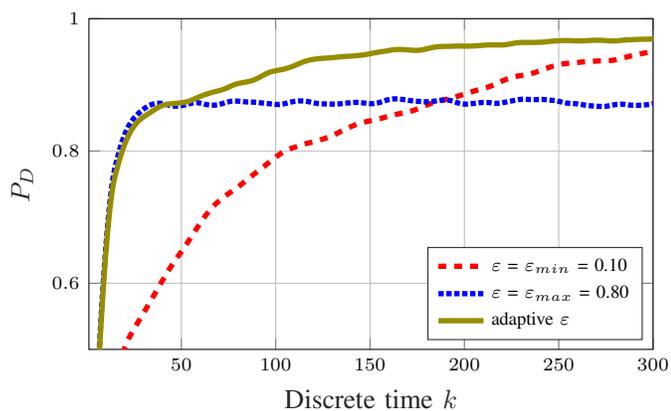Fig. 1: Policy comparison: $P_D$ of target 2 (Scenario 1).



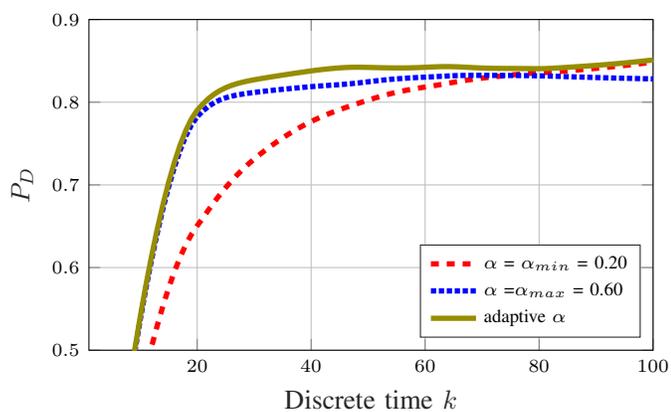Fig. 2: Adaptive vs static $\varepsilon$: $P_D$ of target 2 (Scenario 1).



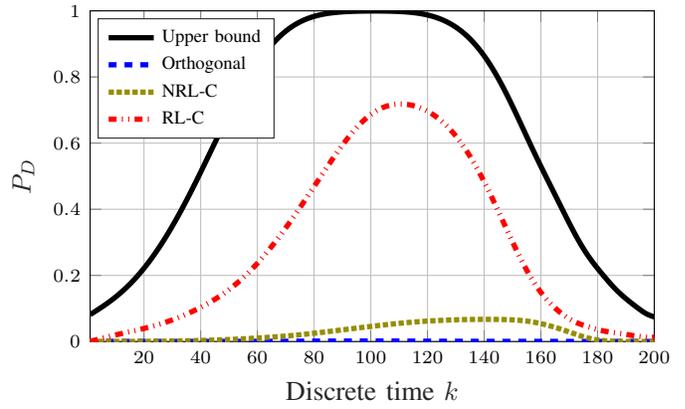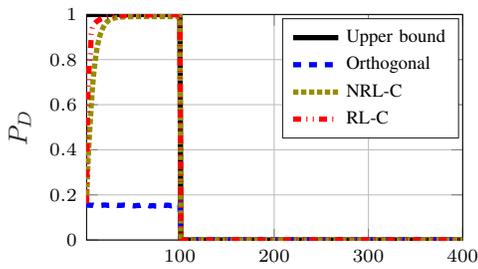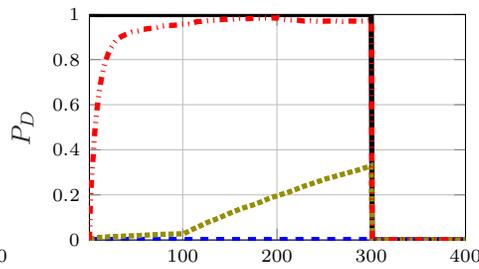Fig. 3: Adaptive vs static $\alpha$: $P_D$ of target 1 (Scenario 2).

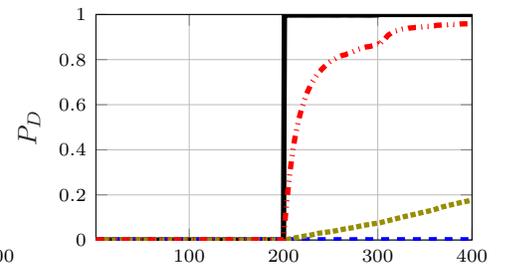Fig. 4: $P_D$ of the targets in Scenario 3. The inset figure in (a) shows the SNR of both targets expressed in dB.



Fig. 5: $P_D$ of the targets in Scenario 4.
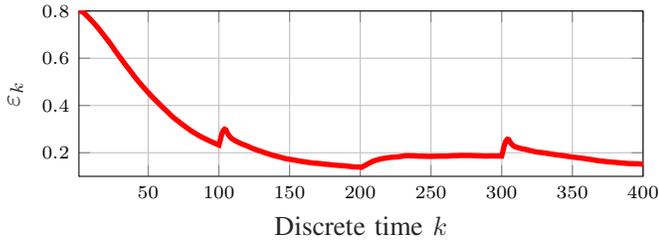


Fig. 6: $\varepsilon_k$ sequence (Scenario 4)



Fig. 7: $\alpha_k$ sequence (Scenario 4)