

## ACTOR-CRITIC ALGORITHMS WITH $\epsilon$ -GREEDY GAUSSIAN POLICY IN MULTIDIMENSIONAL CONTINUOUS ACTION SPACES

CHUNYUAN ZHANG<sup>1,2</sup>, QINGXIN ZHU<sup>1</sup>, YIGUI OU<sup>2</sup> AND XINZHENG NIU<sup>1</sup>

<sup>1</sup>School of Computer Science and Engineering  
University of Electronic Science and Technology of China  
No. 2006, Xiyuan Ave., West Hi-Tech Zone, Chengdu 611731, P. R. China  
zcy7566@126.com; {qxzhu; xinzhengniu}@uestc.edu.cn

<sup>2</sup>College of Information Science and Technology  
Hainan University  
No. 58, Renmin Ave., Meilan District, Haikou 570228, P. R. China  
ouyigui@126.com

Received December 2015; revised March 2016

**ABSTRACT.** *In actor-critic (AC) algorithms, the Gaussian policy is widely used for solving the sequential decision problems with continuous action spaces. However, this policy has a tendency of over-exploration due to lack of greediness, which often makes AC algorithms difficult to obtain good convergence speed and quality. In this paper, we propose a novel  $\epsilon$ -greedy Gaussian policy, and present two compatible AC algorithm frameworks for successfully using it. The proposed policy can be viewed as a hybrid of the traditional Gaussian policy and the  $\epsilon$ -greedy policy. At each time step, it generates some candidate actions by performing symmetric Gaussian perturbations on the current action mean, and then uses the  $\epsilon$ -greedy policy to select the behaviour action based on advantage functions. To the best of our knowledge, this is the first time to introduce the  $\epsilon$ -greedy policy into AC algorithms for solving multidimensional continuous action problems. Theoretical analysis shows that compatible AC algorithms can obtain better convergence quality with the proposed policy than with the traditional Gaussian policy. Finally, experimental results on a mountain-car problem and a puddle-world problem demonstrate the effectiveness of the proposed policy and compatible AC algorithm frameworks.*

**Keywords:** Actor-critic, Gaussian policy, Continuous space, Compatible function approximation, Policy gradient, Reinforcement learning

**1. Introduction.** In the real world, many sequential decision problems have continuous state and action spaces. Currently, there are three main categories of reinforcement learning (RL) methods for solving such problems, that is, actor-only, critic-only and actor-critic (AC) algorithms [1]. Typically, actor-only algorithms work with a parameterized family of policies over which the optimal policy can be searched directly, critic-only algorithms discretize the continuous action space and work with state-action value functions over which the optimal policy can be derived indirectly, and AC algorithms can be viewed as the combination of actor-only and critic-only algorithms. In general, AC algorithms have two separate memory structures, namely, an actor for policy improvement and a critic for policy evaluation [2]. This property makes AC algorithms easy to combine the strong points of actor-only and critic-only algorithms [3]. By introducing the state value function to reduce the variance of policy gradient, AC algorithms have better convergence performance than actor-only algorithms [4]. By representing and improving the policy

explicitly, AC algorithms have more reliable convergence guarantees [3, 4] and higher control precision [5] than critic-only algorithms. Therefore, AC algorithms have become a preferred RL method for continuous action problems.

However, AC algorithms often converge slower and worse than critic-only algorithms in many cases [6, 7]. In recent years, some studies have attempted to overcome this problem. Melo and Lopes proposed a fitted natural AC algorithm by using importance sampling [8], but it is only suitable for batch learning and the implementation is complicated. Wawrzynski and Tanwani proposed two sequential AC algorithms with experience replay [9, 10], whereas they require a large amount of memory to cache history samples. Hasselt and Wiering proposed a continuous AC learning automaton algorithm by using the action error and the sign of the temporal-difference (TD) error for policy update [11, 12]. However, it lacks reliable guarantees in terms of near-optimality of the resulting policy, and it may perform badly in problems with specific types of noise [13]. Vien and Ertel proposed a compatible AC TAMER framework [14], which requires human feedback and thus deviates from the original intention of RL. Lazaric et al. proposed an AC algorithm in which the policy is estimated by sequential Monte Carlo methods [5]. Unfortunately, it is built on a discrete state space, which narrows the applicability in practice. Kimura and Kobayashi proposed an AC algorithm using a binary tree action selector [15], but it does not truly offer continuous actions since the selector is actually a variant of the Boltzmann policy.

Different from the above studies, we argue that the slow convergence of AC algorithms is mainly caused by the inherent limitations of their exploration policies. For continuous action problems, AC algorithms widely use a parameterized normal distribution to represent the exploration policy [16, 17, 18, 19, 20, 21], which is known as the Gaussian policy. In order to distinguish it from our policy proposed in Section 3.3, we call it the traditional Gaussian policy in this paper. At each time step, the traditional Gaussian policy only generates one action by performing a Gaussian perturbation on the current action mean, and has no choice but to use it as the behaviour action. Although we can adaptively adjust the perturbation size by using the state-dependent exploration [22], the Gaussian perturbation still has strong randomness. By contrast, critic-only algorithms widely use an  $\epsilon$ -greedy policy as the exploration policy, which greedily selects the behaviour action with probability  $1 - \epsilon$ . In an  $\epsilon$ -greedy policy,  $\epsilon$  denotes the exploration strength and is often small. Thus, compared with an  $\epsilon$ -greedy policy, the traditional Gaussian policy obviously lacks greediness and has a tendency of over-exploration. In RL, it is well known that over-exploration may result in sample inefficiency [2]. This also explains why AC algorithms often converge slower than critic-only algorithms.

Intuitively, we may improve the convergence performance of AC algorithms by increasing the greediness of the traditional Gaussian policy. However, the traditional Gaussian policy lacks the ability to select the greedy action, and a full search in a continuous action space to find the optimal action is often unfeasible [5, 23]. To deal with these two problems, we proposed a novel AC algorithm framework with symmetric perturbation sampling for one-dimensional continuous action problems in our previous work [24]. At each time step, this framework generates two candidate actions by performing two symmetric Gaussian perturbations on the current action mean, and takes them to interact with the environment in parallel. Then, it selects the behaviour action and updates the value-function parameters based on their maximum TD error, and updates the traditional Gaussian policy parameters based on their average regular gradient or natural gradient. Theoretical analysis and experimental results show that this framework has reliable and fast convergence performance. Nevertheless, it requires two interactions per time step,

which undermines the practical applicability seriously. If we extend it for  $n$ -dimensional continuous action problems, it will require  $2^n$  interactions per time step.

In this paper, we further investigate how to increase the greediness of the traditional Gaussian policy without multi-interactions. The main contributions of this paper include three aspects. (i) We propose a novel  $\epsilon$ -greedy Gaussian policy, which can be viewed as a hybrid of the traditional Gaussian policy and the  $\epsilon$ -greedy policy. (ii) For successfully using this policy, we present two compatible AC algorithm frameworks called RLSEGAC and SEGAC. Compared with our previous framework, they are mainly built on compatible function approximation rather than value function approximation. In particular, they only require one interaction per time step, even for  $n$ -dimensional continuous action problems. (iii) We prove that compatible AC algorithms can obtain better convergence quality with this policy than with the traditional Gaussian policy in theory. The rest of this paper is organized as follows. In Section 2, some background materials are provided. In Section 3, the traditional Gaussian policy and the  $\epsilon$ -greedy policy are reviewed, and then the  $\epsilon$ -greedy Gaussian policy is proposed. In Section 4, two compatible AC algorithm frameworks are presented for using the proposed policy. In Sections 5 and 6, the performance of compatible AC algorithms with the proposed policy is analyzed theoretically and demonstrated experimentally. Finally, conclusions are summarized in Section 7.

**2. Background.** In this section, we introduce the basic definitions, assumptions and notations, which will be used throughout the paper without any further mention. We also review the policy gradient theorem and conventional AC algorithms, which are needed to establish our main results described in Sections 3-5.

**2.1. Preliminaries.** In this paper, we assume that the sequential decision problem with a continuous state space  $\mathcal{S} \in \mathbb{R}^m$  and a continuous action space  $\mathcal{A} \in \mathbb{R}^n$  can be modeled as a discrete-time Markov decision process [16]. For any state  $\mathbf{s}_t \in \mathcal{S}$  at time step  $t \in \mathbb{N}$ , the agent takes an action  $\mathbf{a}_t \in \mathcal{A}$  drawn from a stochastic parameterized policy  $\pi_{\boldsymbol{\theta}}(\mathbf{a}_t|\mathbf{s}_t) = p(\mathbf{a}_t|\mathbf{s}_t, \boldsymbol{\theta})$  with a parameter vector  $\boldsymbol{\theta} \in \mathbb{R}^N$ , then transfers to a successor state  $\mathbf{s}_{t+1}$  drawn from the state transition distribution  $\tau(\mathbf{s}_{t+1}|\mathbf{s}_t, \mathbf{a}_t)$ , and earns a reward  $r_t = r(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \in \mathbb{R}$ . In addition, we assume that  $\pi_{\boldsymbol{\theta}}$  is continuously differentiable w.r.t.  $\boldsymbol{\theta}$ , and use  $\pi$  to denote  $\pi_{\boldsymbol{\theta}}$  for simplicity of notation.

Under the above assumptions, the goal of RL is to learn an optimal  $\boldsymbol{\theta}^*$  for maximizing the long-run expected return  $J^\pi$ . Depending on the nature of the problem, there are two formulations for defining  $J^\pi$  [21, 25]. One is the discount-return formulation, i.e.,  $J^\pi = \mathbf{E}_\pi[\sum_{t=0}^\infty \gamma^t r_t | \mathbf{s}_0]$ , where  $\gamma \in [0, 1]$  is the discount factor and  $\mathbf{s}_0$  is an initial state drawn from the start-state distribution  $p(\mathbf{s}_0)$ . The other is the average-return formulation, i.e.,  $J^\pi = \lim_{l \rightarrow \infty} \frac{1}{l} \mathbf{E}_\pi [\sum_{t=0}^{l-1} r_t]$ . Both formulations can be unifiedly described as

$$J^\pi = \int_{\mathcal{S}} d^\pi(\mathbf{s}) \int_{\mathcal{A}} \pi(\mathbf{a}|\mathbf{s}) R(\mathbf{s}, \mathbf{a}) d\mathbf{a} d\mathbf{s} \tag{1}$$

where  $R(\mathbf{s}, \mathbf{a}) = \mathbf{E}_\tau[r_t | \mathbf{s}_t = \mathbf{s}, \mathbf{a}_t = \mathbf{a}]$  is the expected reward in  $(\mathbf{s}, \mathbf{a})$ , and  $d^\pi(\mathbf{s})$  is the state distribution. For the discount-return formulation,  $d^\pi(\mathbf{s}) = \sum_{t=0}^\infty \gamma^t P(\mathbf{s}_t = \mathbf{s} | \mathbf{s}_0, \pi)$ , whereas, for the average-return formulation,  $d^\pi(\mathbf{s}) = \lim_{t \rightarrow \infty} P(\mathbf{s}_t = \mathbf{s} | \mathbf{s}_0, \pi)$ .

Besides using  $J^\pi$  as the performance metric of  $\pi$ , many RL algorithms also use the state value function  $V^\pi(\mathbf{s})$  or the state-action value function  $Q^\pi(\mathbf{s}, \mathbf{a})$  as the performance metric of  $\pi$  in  $\mathbf{s}$  or  $(\mathbf{s}, \mathbf{a})$ . Similar to  $J^\pi$ ,  $V^\pi(\mathbf{s})$  and  $Q^\pi(\mathbf{s}, \mathbf{a})$  can be defined as follows

$$V^\pi(\mathbf{s}) = \mathbf{E}_\pi \left[ \sum_{t=0}^\infty (\gamma^t r_t - \bar{r}_\pi) \mid \mathbf{s}_0 = \mathbf{s} \right] \tag{2}$$

$$Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbf{E}_\pi \left[ \sum_{t=0}^{\infty} (\gamma^t r_t - \bar{r}_\pi) \mid \mathbf{s}_0 = \mathbf{s}, \mathbf{a}_0 = \mathbf{a} \right] \quad (3)$$

where, for the discount-return formulation,  $\bar{r}_\pi = 0$ , whereas, for the average-return formulation,  $\bar{r}_\pi = J^\pi$  and  $\gamma = 1$ .

**2.2. Policy gradient theorem.** In general,  $\boldsymbol{\theta}^*$  cannot be solved by maximizing  $J^\pi$  directly. Instead, we often resort to policy gradient algorithms, in which  $\boldsymbol{\theta}$  is updated by

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t^\theta \nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi \quad (4)$$

where  $\alpha_t^\theta$  is the step size, and  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$  is the policy gradient.

$\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$  has two main types, namely, the regular gradient  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$  and the natural gradient  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$ . By the policy gradient theorem [25],  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$  can be defined as

$$\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi = \int_{\mathcal{S}} d^\pi(\mathbf{s}) \int_{\mathcal{A}} Q^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s}) d\mathbf{a} d\mathbf{s} = \mathbf{E}_\pi [Q^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s})] \quad (5)$$

Since  $\int_{\mathcal{A}} \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s}) d\mathbf{a} = 0$ , we can introduce an arbitrary baseline  $b(\mathbf{s})$  into  $Q^\pi(\mathbf{s}, \mathbf{a})$  without affecting the unbiasedness of the gradient estimate [4, 16]. Let  $A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - b(\mathbf{s})$ . Then, Equation (5) can be rewritten as

$$\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi = \mathbf{E}_\pi [A^\pi(\mathbf{s}, \mathbf{a}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s})] \quad (6)$$

By linearly transforming  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$ ,  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$  can be defined as [4, 16]

$$\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi = \mathbf{F}^{-1}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi \quad (7)$$

where  $\mathbf{F}(\boldsymbol{\theta}) = \mathbf{E}_\pi [\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s}) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s})^\top]$  is the Fisher information matrix.

In the implementation of the policy gradient theorem, there are two ways to deal with  $A^\pi(\mathbf{s}, \mathbf{a})$  for estimating  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$ . The first way is to approximate  $A^\pi(\mathbf{s}, \mathbf{a})$  with the compatible function [4, 16], and is mainly used for estimating  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$ . Let  $\hat{A}(\mathbf{s}, \mathbf{a}) = \hat{Q}(\mathbf{s}, \mathbf{a}) - \hat{b}(\mathbf{s}) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s})$ , where  $\hat{A}(\mathbf{s}, \mathbf{a})$ ,  $\hat{Q}(\mathbf{s}, \mathbf{a})$  and  $\hat{b}(\mathbf{s})$  denote the approximations of  $A^\pi(\mathbf{s}, \mathbf{a})$ ,  $Q^\pi(\mathbf{s}, \mathbf{a})$  and  $b(\mathbf{s})$ ,  $\mathbf{w} \in \mathbb{R}^N$  denotes the compatible parameter vector, and  $\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s})$  is called the compatible feature vector. Since  $\int_{\mathcal{A}} \pi(\mathbf{a} \mid \mathbf{s}) \hat{A}(\mathbf{s}, \mathbf{a}) d\mathbf{a} = \mathbf{w}^\top \int_{\mathcal{A}} \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s}) d\mathbf{a} = 0$ , we have  $\int_{\mathcal{A}} \pi(\mathbf{a} \mid \mathbf{s}) \hat{Q}(\mathbf{s}, \mathbf{a}) d\mathbf{a} = \int_{\mathcal{A}} \pi(\mathbf{a} \mid \mathbf{s}) \hat{b}(\mathbf{s}) d\mathbf{a} = \hat{b}(\mathbf{s})$ . Further, from Equations (2) and (3), we have  $\int_{\mathcal{A}} \pi(\mathbf{a} \mid \mathbf{s}) Q^\pi(\mathbf{s}, \mathbf{a}) d\mathbf{a} = V^\pi(\mathbf{s})$ . Thus,  $\hat{b}(\mathbf{s})$  is the approximation of  $V^\pi(\mathbf{s})$  [26], and we can assume  $b(\mathbf{s}) = V^\pi(\mathbf{s})$ . In this paper,  $A^\pi(\mathbf{s}, \mathbf{a}) = Q^\pi(\mathbf{s}, \mathbf{a}) - V^\pi(\mathbf{s})$  is the advantage function, which can be approximated by

$$\hat{A}(\mathbf{s}, \mathbf{a}) = \hat{Q}(\mathbf{s}, \mathbf{a}) - \hat{V}(\mathbf{s}) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} \mid \mathbf{s}) \quad (8)$$

where  $\hat{V}(\mathbf{s})$  denotes the approximation of  $V^\pi(\mathbf{s})$ . Then,  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$  can be estimated as

$$\hat{\nabla}_{\boldsymbol{\theta}}^{\text{NG}} J^\pi = \mathbf{w} \quad (9)$$

The second way is to replace  $A^\pi(\mathbf{s}, \mathbf{a})$  with some form of TD errors [4, 21], and thus the estimate of  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$  is associated with  $V^\pi(\mathbf{s})$  and  $\bar{r}_\pi$  rather than  $A^\pi(\mathbf{s}, \mathbf{a})$ . Different from the first way, this way is mainly used for estimating  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$ .

**2.3. Conventional AC algorithms.** As stated in Section 1, an AC algorithm often has two separate memory structures, which make it easy to implement the policy gradient theorem. According to the gradient type, AC algorithms can be divided into regular-gradient AC (RAC) algorithms and natural-gradient AC (NAC) algorithms.

In conventional RAC algorithms, the actor often uses the second way to estimate  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$  for updating  $\boldsymbol{\theta}$ , and the critic often uses linear TD( $\lambda$ ) algorithms [27, 28] for learning  $V^\pi(\mathbf{s})$  and  $\bar{r}_\pi$ . Define the one-step TD error as  $\delta_t^V = r_t - \bar{r}_{t+1} + \gamma \hat{V}(\mathbf{s}_{t+1}) - \hat{V}(\mathbf{s}_t)$ , where  $\bar{r}_{t+1}$

denotes the estimate of  $\bar{r}_\pi$  at time step  $t$ . Let  $\hat{V}(\mathbf{s}) = \mathbf{v}^\top \boldsymbol{\phi}(\mathbf{s})$ , where  $\mathbf{v} \in \mathbb{R}^M$  denotes the parameter vector of  $\hat{V}(\mathbf{s})$ , and  $\boldsymbol{\phi}(\mathbf{s}) \in \mathbb{R}^M$  denotes the feature vector of  $\hat{V}(\mathbf{s})$ . Then,  $\delta_t^V$  can be estimated as

$$\delta_t^V = r_t - \bar{r}_{t+1} + \mathbf{v}_t^\top (\gamma \boldsymbol{\phi}(\mathbf{s}_{t+1}) - \boldsymbol{\phi}(\mathbf{s}_t)) \tag{10}$$

All update rules in the critic are [21]

$$\bar{r}_{t+1} = (1 - \alpha_t^{\bar{r}}) \bar{r}_t + \alpha_t^{\bar{r}} r_t \tag{11}$$

$$\mathbf{z}_{t+1}^v = \gamma \lambda \mathbf{z}_t^v + \boldsymbol{\phi}(\mathbf{s}_t) \tag{12}$$

$$\mathbf{v}_{t+1} = \mathbf{v}_t + \alpha_t^v \delta_t^V \mathbf{z}_{t+1}^v \tag{13}$$

where  $\lambda \in [0, 1]$ ,  $\mathbf{z}_{t+1}^v \in \mathbb{R}^M$ ,  $\alpha_t^{\bar{r}}$  and  $\alpha_t^v$  are the trace-decay parameter, the trace of  $\boldsymbol{\phi}(\mathbf{s}_t)$  and the step sizes, respectively. Let  $\mathbf{z}_{t+1}^\theta$  denote the trace of  $\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t)$ , i.e.,

$$\mathbf{z}_{t+1}^\theta = \gamma \lambda \mathbf{z}_t^\theta + \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \tag{14}$$

Referring to A.2 in [18], we can verify that  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$  can be approximated by

$$\hat{\nabla}_{\boldsymbol{\theta}}^{\text{RG}} J^\pi = \delta_t^V \mathbf{z}_{t+1}^\theta \tag{15}$$

Thus, the update rule in the actor is [21]

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t^\theta \delta_t^V \mathbf{z}_{t+1}^\theta \tag{16}$$

In conventional NAC algorithms, the actor often uses the first way to estimate  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$  for updating  $\boldsymbol{\theta}$ , and the critic often uses Equation (8) to approximate  $A^\pi(\mathbf{s}, \mathbf{a})$ . However,  $A^\pi(\mathbf{s}, \mathbf{a})$  cannot be learned with TD-like bootstrapping exclusively. To remedy this situation, Peters and Schaal proposed a least-squares NAC (LSNAC) algorithm [16], in which the critic uses LSTD-Q( $\lambda$ ) for learning  $A^\pi(\mathbf{s}, \mathbf{a})$ . However, this algorithm has  $O((M+N)^3)$  computational complexity per time step and only considers the discount-return formulation. On this basis, Wawrzynski proposed a sequential NAC (SNAC) algorithm [9], and Degris et al. proposed an incremental NAC (INAC) algorithm [21]. Compared with SNAC, INAC considers both return formulations defined in Section 2.1. Here we only review INAC briefly. The other two algorithms will be left in Section 4 for further discussion. In the critic of INAC, the update rule for  $\mathbf{w}$  is

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t^w (\delta_t^V \mathbf{z}_{t+1}^\theta - \mathbf{F}_t \mathbf{w}_t) \tag{17}$$

where  $\mathbf{F}_t = \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t)^\top$ , and the estimate of  $\delta_t^V$  is the same as Equation (10). The other update rules for  $\bar{r}_t$ ,  $\mathbf{z}_t^v$  and  $\mathbf{v}_t$  are the same as Equations (11)-(13). In the actor of INAC, the update rule is

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t^\theta \mathbf{w}_{t+1} \tag{18}$$

### 3. Exploration Policies.

In this section, we review the traditional Gaussian policy and the  $\epsilon$ -greedy policy, and then propose the  $\epsilon$ -greedy Gaussian policy by combining them.

**3.1. Traditional Gaussian policy.** The traditional Gaussian policy is a parameterized normal distribution, which is widely used in AC algorithms to explore the continuous action space. We assume that all components of the continuous action are independent of each other. Then, an  $n$ -dimensional traditional Gaussian policy can be defined as [22]

$$\pi(\mathbf{a} | \mathbf{s}) = \prod_{i=1}^n \mathcal{N}(\mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu), \sigma_i^2(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)) \tag{19}$$

where  $\mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu)$  and  $\sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$  denote the mean and standard deviation of the  $i$ -th action component  $a_i$  in  $\mathbf{s}$ , and they are parameterized by  $\boldsymbol{\theta}_i^\mu$  and  $\boldsymbol{\theta}_i^\sigma$  respectively. From the above equation,  $a_i$  can be generated as

$$a_i = \mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu) + \rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma) \quad (20)$$

where  $\rho_i \sim \mathcal{N}(0, 1)$ , and  $\rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$  is called the Gaussian perturbation. Correspondingly, the compatible feature vector can be calculated as

$$\begin{cases} \nabla_{\boldsymbol{\theta}_i^\mu} \log \pi(\mathbf{a}|\mathbf{s}) = \frac{a_i \mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu)}{\sigma_i^2(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)} \nabla_{\boldsymbol{\theta}_i^\mu} \mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu) \\ \nabla_{\boldsymbol{\theta}_i^\sigma} \log \pi(\mathbf{a}|\mathbf{s}) = \frac{\rho_i^2 - 1}{\sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)} \nabla_{\boldsymbol{\theta}_i^\sigma} \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma) \end{cases} \quad (21)$$

From Equations (20), (21), (6) and (7),  $\rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$ ,  $a_i$  and  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$  are associated with  $\rho_i$ . Obviously, they have strong randomness since  $\rho_i$  is drawn from  $\mathcal{N}(0, 1)$ . This may result in over-exploration and sample inefficiency. In Section 1 and our previous work [24], we have analyzed the inherent limitations of the traditional Gaussian policy in detail, so we will not repeat them here.

**3.2.  $\epsilon$ -greedy policy.** In RL, the  $\epsilon$ -greedy policy may be the most popular exploration policy, which is widely used in critic-only algorithms to explore the discrete action space. Let  $\mathbf{a}^* = \arg \max_{\mathbf{a} \in \mathcal{A}(\mathbf{s})} Q^\pi(\mathbf{s}, \mathbf{a})$ , where  $\mathcal{A}(\mathbf{s})$  is the discrete action set in  $\mathbf{s}$ . Then, an  $\epsilon$ -greedy policy can be defined as [2]

$$\pi(\mathbf{a}|\mathbf{s}) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(\mathbf{s})|} & \text{if } \mathbf{a} = \mathbf{a}^* \\ \frac{\epsilon}{|\mathcal{A}(\mathbf{s})|} & \text{if } \mathbf{a} \neq \mathbf{a}^* \end{cases} \quad (22)$$

where  $\epsilon \in (0, 1]$  is the exploration factor and denotes the exploration strength.

The  $\epsilon$ -greedy policy has some good properties. For any  $\epsilon$ -soft policy  $\zeta$ , any  $\epsilon$ -greedy policy with respect to  $Q^\zeta$  is guaranteed to be better than or equal to  $\zeta$  [2]. In addition, by selecting a moderate  $\epsilon$ , the  $\epsilon$ -greedy policy can balance exploration and exploitation well. However, the  $\epsilon$ -greedy policy is unsuitable for continuous action problems, since it requires a global maximization at each time step [5, 23].

**3.3.  $\epsilon$ -greedy Gaussian policy.** Since the normal distribution is a symmetric distribution, Equation (20) can also be rewritten as  $a_i = \mu_i(\mathbf{s}, \boldsymbol{\theta}_i^\mu) - \rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$ . In this paper and our previous work,  $\rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$  and  $-\rho_i \sigma_i(\mathbf{s}, \boldsymbol{\theta}_i^\sigma)$  are called symmetric Gaussian perturbations. Thus, an  $n$ -dimensional traditional Gaussian policy can be equivalently implemented as follows. We first generate  $2^n$  candidate actions by using symmetric Gaussian perturbations as shown in Figure 1, and then randomly select a candidate as the behaviour action. Intuitively, if we can greedily select the behaviour action from candidates, we will reduce the randomness of the traditional Gaussian policy. In our previous work [24], we proposed an AC algorithm framework following this idea, and demonstrated its effectiveness in the one-dimensional action space. However, as mentioned in Section 1, this framework requires  $2^n$  interactions per time step for  $n$ -dimensional continuous action problems, since it selects the behaviour action based on TD errors. From Equation (8), we have  $\arg \max_{\mathbf{a} \in \mathcal{A}} \hat{Q}(\mathbf{s}, \mathbf{a}) = \arg \max_{\mathbf{a} \in \mathcal{A}} \hat{A}(\mathbf{s}, \mathbf{a})$ . If we use advantage functions to replace TD errors for selecting the behaviour action, we will avoid this situation.

Based on the above analysis, we propose a novel exploration policy, called the  $\epsilon$ -greedy Gaussian policy, for exploring the  $n$ -dimensional continuous action space. It can be viewed

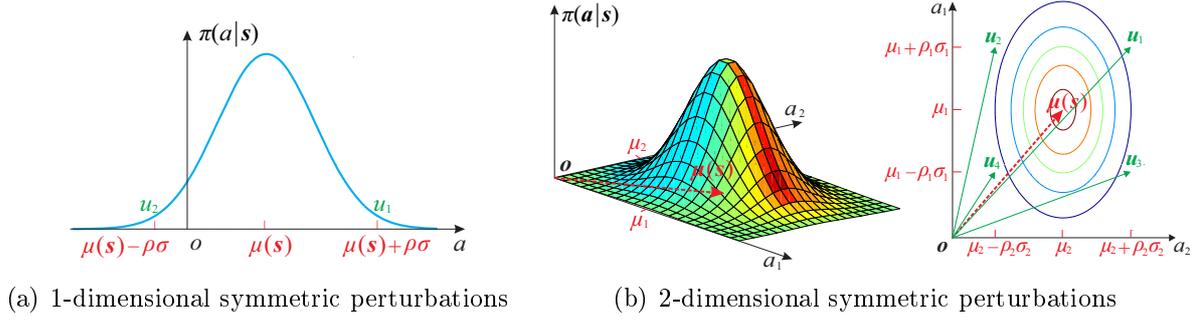


FIGURE 1. 1-dimensional and 2-dimensional symmetric Gaussian perturbations

as a hybrid of the traditional Gaussian policy and the  $\epsilon$ -greedy policy. Firstly, by performing symmetric Gaussian perturbations on the current action mean, it generates a candidate action set  $\mathcal{A}(\mathbf{s}) = \{\mathbf{u}_1, \dots, \mathbf{u}_{2^n}\}$ , where the  $j$ -th candidate action is defined as

$$\mathbf{u}_j = \boldsymbol{\mu}(\mathbf{s}) + \left[ (-1)^{\lfloor \frac{j-1}{2^{n-1}} \rfloor}, \dots, (-1)^{\lfloor \frac{j-1}{2^{n-n}} \rfloor} \right]^\top \circ \boldsymbol{\rho} \circ \boldsymbol{\sigma}(\mathbf{s}) \quad (23)$$

where  $\boldsymbol{\mu}(\mathbf{s}) = [\mu_1(\mathbf{s}, \boldsymbol{\theta}_1^\mu), \dots, \mu_n(\mathbf{s}, \boldsymbol{\theta}_n^\mu)]^\top$  and  $\boldsymbol{\sigma}(\mathbf{s}) = [\sigma_1(\mathbf{s}, \boldsymbol{\theta}_1^\sigma), \dots, \sigma_n(\mathbf{s}, \boldsymbol{\theta}_n^\sigma)]^\top$  denote the mean and standard deviation of the current action,  $\boldsymbol{\rho}$  is drawn from the  $n$ -dimensional standard normal distribution  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ , and  $\circ$  denotes the Hadamard product. Secondly, it selects the behaviour action  $\mathbf{a}$  from  $\mathcal{A}(\mathbf{s})$  with the  $\epsilon$ -greedy policy, where  $\mathbf{a}^* = \arg \max_{\mathbf{u}_j \in \mathcal{A}(\mathbf{s})} \hat{A}(\mathbf{s}, \mathbf{u}_j)$ . The complete algorithm of this policy is shown in Algorithm 1.

---

**Algorithm 1**  $\epsilon$ -greedy Gaussian policy
 

---

- 1: **Input:**  $\epsilon, \mathbf{s}, \mathbf{w}, \boldsymbol{\theta}, \phi(\mathbf{s}), \pi(\mathbf{u}|\mathbf{s}), \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}|\mathbf{s})$
  - 2: **Initialize:**  $\mathcal{A}(\mathbf{s}) = \{\}$
  - 3: Approximate  $\boldsymbol{\mu}(\mathbf{s})$  and  $\boldsymbol{\sigma}(\mathbf{s})$  with  $\boldsymbol{\theta}$  and  $\phi(\mathbf{s})$
  - 4: Draw  $\boldsymbol{\rho} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: **for**  $j = 1, 2, \dots, 2^n$  **do**
  - 6:     Generate candidate action  $\mathbf{u}_j$  by Equation (23)
  - 7:      $\hat{A}(\mathbf{s}, \mathbf{u}_j) = \mathbf{w}^\top \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{u}_j|\mathbf{s})$
  - 8:      $\mathcal{A}(\mathbf{s}) = \mathcal{A}(\mathbf{s}) \cup \{\mathbf{u}_j\}$
  - 9: **end for**
  - 10: Select behaviour action  $\mathbf{a}$  from  $\mathcal{A}(\mathbf{s})$  with the  $\epsilon$ -greedy policy
- 

For the sake of discussion, let  $\vartheta$ ,  $\kappa$  and  $\pi$  denote the traditional Gaussian policy, the  $\epsilon$ -greedy policy and the  $\epsilon$ -greedy Gaussian policy in the rest of paper. Based on the above description, there are  $2^n$  candidate actions in  $\mathbf{s}$ , and their probability densities are  $\vartheta(\mathbf{a}|\mathbf{s})$ . Thus,  $\pi(\mathbf{a}|\mathbf{s})$  can be represented as

$$\pi(\mathbf{a}|\mathbf{s}) = 2^n \vartheta(\mathbf{a}|\mathbf{s}) \kappa(\mathbf{a}|\mathbf{s}) \quad (24)$$

From Equations (23) and (24), the  $\epsilon$ -greedy Gaussian policy has the following properties. (i) It has the same compatible feature vector as the traditional Gaussian policy, i.e.,  $\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s}) = \nabla_{\boldsymbol{\theta}} \log \vartheta(\mathbf{a}|\mathbf{s})$ , since  $\kappa(\mathbf{a}|\mathbf{s})$  is independent of  $\boldsymbol{\theta}$ . (ii) If  $\epsilon = 1$ , it will be degenerated into the traditional Gaussian policy. (iii) All candidate actions have the same generating probability  $\vartheta(\mathbf{a}|\mathbf{s})$  and perturbation size  $|\boldsymbol{\rho} \circ \boldsymbol{\sigma}(\mathbf{s})|$ . (iv) The perturbation directions of  $\mathbf{u}_j$  and  $\mathbf{u}_{2^n-j+1}$  are completely opposite. (v)  $\pi(\mathbf{a}^*|\mathbf{s}) = 2^n \left(1 - \epsilon + \frac{\epsilon}{2^n}\right) \vartheta(\mathbf{a}|\mathbf{s})$ .

**4. Compatible AC Algorithms.** To successfully use the  $\epsilon$ -greedy Gaussian policy for multidimensional continuous action problems, it is necessary to address the learning of the advantage function. In Section 2.3, we have reviewed three NAC algorithms for tackling this problem. In theory, the  $\epsilon$ -greedy Gaussian policy can be used in them seamlessly. However, they only consider the natural gradient, and most of them only consider one return formulation. Although INAC considers both formulations, the quality of the resulting policy may be significantly affected by the choice of  $\phi(\mathbf{s})$  [29]. In our previous work [24] and in [21], we and Degris et al. also find that INAC does not perform very well. In this section, we present two general AC frameworks with the  $\epsilon$ -greedy Gaussian policy by considering compatible function approximation. The first framework is called RLSEGAC, which is built by extending the LSNAC algorithm [16]. The second framework is called SEGAC, which is built by extending the SNAC algorithm [9].

**4.1. RLSEGAC framework.** Now we introduce our first framework. Different from the derivation of LSNAC, our derivation is based on Lemma 3 in [4] rather than the Bellman equation. In addition, our framework uses the recursive least-squares (RLS) method [30], and thus has  $O((M+N)^2)$  computational complexity per time step.

From Lemma 3 in [4], we have  $A^\pi(\mathbf{s}_t, \mathbf{a}_t) = \mathbf{E}_\tau [\delta_t^V | \mathbf{s}_t, \mathbf{a}_t, \pi]$ . Then, Equation (8) can be rewritten as

$$\mathbf{w}_t^\top \nabla_\theta \log \pi(\mathbf{a}_t | \mathbf{s}_t) = \delta_t^V + \eta(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \quad (25)$$

where  $\mathbf{E}_\tau [\eta(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1})] = 0$ . Let  $\hat{V}(\mathbf{s}) = \mathbf{v}^\top \phi(\mathbf{s})$ ,  $\boldsymbol{\psi}_t = [\phi(\mathbf{s}_t)^\top, \nabla_\theta \log \pi(\mathbf{a}_t | \mathbf{s}_t)^\top]^\top$ ,  $\boldsymbol{\varphi}_t = [\phi(\mathbf{s}_{t+1})^\top, \mathbf{0}^\top]^\top$ ,  $\boldsymbol{\xi}_t = [\mathbf{v}_t^\top, \mathbf{w}_t^\top]^\top$  and  $\mathbf{z}_{t+1}^\xi = [\mathbf{z}_{t+1}^v{}^\top, \mathbf{z}_{t+1}^\theta{}^\top]^\top = \gamma \lambda \mathbf{z}_t^\xi + \boldsymbol{\psi}_t$ . Plugging Equation (10) into Equation (25), we can get

$$(\boldsymbol{\psi}_t - \gamma \boldsymbol{\varphi}_t)^\top \boldsymbol{\xi}_t = r_t - \bar{r}_{t+1} + \eta(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \quad (26)$$

Multiplying both sides by  $\mathbf{z}_{t+1}^\xi$ , we have

$$\mathbf{z}_{t+1}^\xi (\boldsymbol{\psi}_t - \gamma \boldsymbol{\varphi}_t)^\top \boldsymbol{\xi}_t = \mathbf{z}_{t+1}^\xi (r_t - \bar{r}_{t+1}) + \eta(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \mathbf{z}_{t+1}^\xi \quad (27)$$

Since  $\mathbf{z}_{t+1}^\xi$  is independent of the state transition distribution  $\tau(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ , we have  $\mathbf{E}_\tau [\eta(\mathbf{s}_t, \mathbf{a}_t, \mathbf{s}_{t+1}) \mathbf{z}_{t+1}^\xi] = \mathbf{0}$ . Then, we can obtain

$$\mathbf{E}_\tau [\mathbf{z}_{t+1}^\xi (\boldsymbol{\psi}_t - \gamma \boldsymbol{\varphi}_t)^\top \boldsymbol{\xi}_t] = \mathbf{E}_\tau [\mathbf{z}_{t+1}^\xi (r_t - \bar{r}_{t+1})] \quad (28)$$

Thus, we can define the optimal function for  $\boldsymbol{\xi}_t$  as  $\varepsilon_t^\pi(\boldsymbol{\xi}) = \|\mathbf{A}_t \boldsymbol{\xi} - \mathbf{b}_t\|_2^2$ , where  $\mathbf{A}_t = \sum_{i=0}^t \mathbf{z}_{i+1}^\xi (\boldsymbol{\psi}_i - \gamma \boldsymbol{\varphi}_i)^\top$  and  $\mathbf{b}_t = \sum_{i=0}^t \mathbf{z}_{i+1}^\xi (r_i - \bar{r}_{i+1})$ . By minimizing  $\varepsilon_t^\pi(\boldsymbol{\xi})$ , we can get

$$\boldsymbol{\xi}_t = \mathbf{A}_t^{-1} \mathbf{b}_t \quad (29)$$

To reduce the computational burden of matrix inversion, we use the RLS method to solve Equation (29). Let  $\mathbf{P}_t = \mathbf{A}_t^{-1}$ ,  $\mathbf{P}_0 = \zeta \mathbf{E}$ ,  $\mathbf{h}_t = (\boldsymbol{\psi}_t - \gamma \boldsymbol{\varphi}_t)^\top \mathbf{P}_t$ ,  $\mathbf{g}_t = \mathbf{P}_t \mathbf{z}_{t+1}^\xi$  and  $k_t = 1 + \mathbf{h}_t \mathbf{z}_{t+1}^\xi$ , where  $\zeta$  is a positive number and  $\mathbf{E}$  is the identity matrix. Then, the update rules for  $\mathbf{P}_t$  and  $\boldsymbol{\xi}_t$  in the critic of RLSEGAC are

$$\mathbf{P}_{t+1} = \mathbf{P}_t - \frac{1}{k_t} \mathbf{g}_t \mathbf{h}_t \quad (30)$$

$$\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t + \frac{1}{k_t} \delta_t^Q \mathbf{g}_t \quad (31)$$

where  $\delta_t^Q$  is defined as

$$\delta_t^Q = r_t - \bar{r}_{t+1} + \boldsymbol{\xi}_t^\top (\gamma \boldsymbol{\varphi}_t - \boldsymbol{\psi}_t) \quad (32)$$

From Equation (8), we have  $\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = \boldsymbol{\xi}_t^\top \boldsymbol{\psi}_t$ . Thus, Equation (32) can be rewritten as

$$\delta_t^Q = \delta_t^V + \hat{V}(\mathbf{s}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = \delta_t^V - \hat{A}(\mathbf{s}_t, \mathbf{a}_t) \quad (33)$$

Since we use compatible function approximation in the critic,  $\nabla_{\boldsymbol{\theta}}^{\text{NG}} J^\pi$  can be estimated by Equation (9). Then, from Equation (7),  $\nabla_{\boldsymbol{\theta}}^{\text{RG}} J^\pi$  can be estimated by  $\hat{\nabla}_{\boldsymbol{\theta}}^{\text{RG}} J^\pi = \mathbf{F}_t \mathbf{w}_{t+1}$  at time step  $t$ . For further reducing the computational cost, we rewrite  $\mathbf{F}_t \mathbf{w}_{t+1} = \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t)^\top \mathbf{w}_{t+1} = \mathbf{w}_{t+1}^\top \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t)$ . Thus, the regular-gradient update rule for  $\boldsymbol{\theta}$  is

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha_t^\theta \mathbf{w}_{t+1}^\top \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}_t | \mathbf{s}_t) \quad (34)$$

and the natural-gradient update rule for  $\boldsymbol{\theta}$  is the same as Equation (18).

It is important to note that here the critic is for a fixed policy since it approximates  $Q^\pi(\mathbf{s}, \mathbf{a})$  rather than  $V^\pi(\mathbf{s})$ . That means samples generated by previous policies are not very good for improving subsequent policies [16, 31, 32, 33]. Comparing this critic with the critics of conventional RAC algorithms and the INAC algorithm, we can find that this adverse effect is mainly caused by the introduction of  $\hat{A}(\mathbf{s}, \mathbf{a})$ . To overcome this problem, Morimura et al. [33] and Wawrzynski [9] suggested to decay  $\mathbf{w}$  after each update, i.e.,

$$\mathbf{w}_{t+1} = \beta \mathbf{w}_{t+1} \quad (35)$$

where  $\beta \in [0, 1]$  is the forgetting factor. Based on the above derivation, the complete RLSEGAC framework is shown in Algorithm 2.

---

**Algorithm 2** RLSEGAC with the  $\epsilon$ -greedy Gaussian policy
 

---

- 1: **Input:**  $\phi(\mathbf{s}), \pi(\mathbf{a} | \mathbf{s}), \nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a} | \mathbf{s})$
  - 2: **Initialize:**  $\epsilon, \gamma, \lambda, \beta, \alpha_0^\pi, \alpha_0^\theta, \mathbf{z}_0^v, \mathbf{z}_0^\theta, \bar{r}_0, \mathbf{v}_0, \mathbf{w}_0, \boldsymbol{\theta}_0, \mathbf{P}_0, \mathbf{s}_0$
  - 3: **for**  $t = 0, 1, 2, \dots$  **do**
  - 4:     Draw action  $\mathbf{a}_t$  by Algorithm 1
  - 5:     Take  $\mathbf{a}_t$ , observe next state  $\mathbf{s}_{t+1}$  and reward  $r_t$
  - 6:     Update  $\bar{r}_t, \mathbf{z}_t^v$  and  $\mathbf{z}_t^\theta$  by Equations (11), (12) and (14)
  - 7:     Compute  $\delta_t^Q$  by Equation (32)
  - 8:     Compute  $\mathbf{h}_t, \mathbf{g}_t$  and  $k_t$
  - 9:     Update  $\mathbf{P}_t, \mathbf{v}_t$  and  $\mathbf{w}_t$  by Equations (30) and (31)
  - 10:    Update  $\boldsymbol{\theta}_t$  by Equation (34) or (18)
  - 11:    Decay  $\mathbf{w}_{t+1}$  by (35)
  - 12:    Update  $\alpha_t^\pi$  and  $\alpha_t^\theta$
  - 13: **end for**
- 

**4.2. SEGAC framework.** In this subsection, we present our second framework. It can be viewed as an incremental version of RLSEGAC, and has  $O(M + N)$  computational complexity per time step.

Let  $\boldsymbol{\xi}_t, \boldsymbol{\psi}_t, \mathbf{z}_t^\xi$  and  $\delta_t^Q$  have the same meaning as those in Section 4.1. Defining the optimal function  $\varepsilon^\pi(\boldsymbol{\xi}_t) = \mathbf{E}_\pi \left[ \left( \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) - Q^\pi(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \middle| \mathbf{s}_t, \mathbf{a}_t \right]$ , we have

$$\nabla_{\boldsymbol{\xi}_t} \varepsilon^\pi(\boldsymbol{\xi}_t) = 2 \mathbf{E}_\pi \left[ \left( \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) - Q^\pi(\mathbf{s}_t, \mathbf{a}_t) \right) \boldsymbol{\psi}_t \middle| \mathbf{s}_t, \mathbf{a}_t \right] \quad (36)$$

Furthermore, define the  $\lambda$ -return as

$$R_t^\lambda = \sum_{n=1}^{\infty} (1 - \lambda) \lambda^{n-1} R_t^{(n)} \quad (37)$$

where  $R_t^{(n)}$  denotes the  $n$ -step return [2], which is defined as

$$R_t^{(n)} = \begin{cases} r_t - \bar{r}_{t+1} + \gamma \hat{V}(\mathbf{s}_{t+1}) & \text{if } n = 1 \\ r_t - \bar{r}_{t+1} + \gamma \hat{V}(\mathbf{s}_{t+1}) + \sum_{i=1}^{n-1} \gamma^i \delta_{t+i+1}^Q & \text{if } n > 1 \end{cases}$$

Let  $\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = \boldsymbol{\xi}_t^\top \boldsymbol{\psi}_t$ ,  $\hat{V}(\mathbf{s}) = \mathbf{v}^\top \boldsymbol{\phi}(\mathbf{s})$  and  $Q^\pi(\mathbf{s}_t, \mathbf{a}_t) \approx R_t^\lambda$ . Then, by the forward-backward view analysis [2, 18],  $\nabla_{\boldsymbol{\xi}_t} \varepsilon^\pi(\boldsymbol{\xi}_t)$  can be estimated as

$$\hat{\nabla}_{\boldsymbol{\xi}_t} \varepsilon^\pi(\boldsymbol{\xi}_t) = -2\delta_t^Q \mathbf{z}_{t+1}^\xi \quad (38)$$

Thus, the update rule for  $\boldsymbol{\xi}$  in the critic of SEGAC is

$$\boldsymbol{\xi}_{t+1} = \boldsymbol{\xi}_t + \alpha_t^\xi \delta_t^Q \mathbf{z}_{t+1}^\xi \quad (39)$$

where  $\alpha_t^\xi$  is the step size.

The actor of SEGAC is the same as that of RLSEGAC. For the sake of brevity, we omit the derivation. Finally, we summarize the complete SEGAC framework in Algorithm 3.

---

**Algorithm 3** SEGAC with the  $\epsilon$ -greedy Gaussian policy

---

- 1: **Input:**  $\boldsymbol{\phi}(\mathbf{s})$ ,  $\pi(\mathbf{a}|\mathbf{s})$ ,  $\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s})$
  - 2: **Initialize:**  $\epsilon$ ,  $\gamma$ ,  $\lambda$ ,  $\beta$ ,  $\alpha_0^{\bar{r}}$ ,  $\alpha_0^\xi$ ,  $\alpha_0^\theta$ ,  $\mathbf{z}_0^v$ ,  $\mathbf{z}_0^\theta$ ,  $\bar{r}_0$ ,  $\mathbf{v}_0$ ,  $\mathbf{w}_0$ ,  $\boldsymbol{\theta}_0$ ,  $\mathbf{s}_0$
  - 3: **for**  $t = 0, 1, 2, \dots$  **do**
  - 4:   Draw action  $\mathbf{a}_t$  by Algorithm 1
  - 5:   Take  $\mathbf{a}_t$ , observe next state  $\mathbf{s}_{t+1}$  and reward  $r_t$
  - 6:   Update  $\bar{r}_t$ ,  $\mathbf{z}_t^v$  and  $\mathbf{z}_t^\theta$  by Equations (11), (12) and (14)
  - 7:   Compute  $\delta_t^Q$  by Equation (32)
  - 8:   Update  $\mathbf{v}_t$  and  $\mathbf{w}_t$  by Equation (39)
  - 9:   Update  $\boldsymbol{\theta}_t$  by Equation (34) or (18)
  - 10:   Decay  $\mathbf{w}_{t+1}$  by (35)
  - 11:   Update  $\alpha_t^{\bar{r}}$ ,  $\alpha_t^\xi$  and  $\alpha_t^\theta$
  - 12: **end for**
- 

**Remark 4.1.** For episodic problems, when  $\mathbf{s}_{t+1}$  is an absorbing state, RLSEGAC and SEGAC should let  $\boldsymbol{\varphi}_t = \mathbf{0}$ , and reset  $\mathbf{s}_{t+1}$  as the start state of next episode after  $\boldsymbol{\theta}_t$  is updated. For the average-return formulation, there is no  $\bar{r}_t$  and no  $\alpha_t^{\bar{r}}$  in RLSEGAC and SEGAC since  $\bar{r}_\pi = 0$ .

**Remark 4.2.** If  $\epsilon = 1$  or  $\mathbf{a}_t$  is drawn from the traditional Gaussian policy, RLSEGAC and SEGAC will be degenerated into AC algorithm frameworks with the traditional Gaussian policy, which are called RLSTAC and STAC in this paper.

**Remark 4.3.** Different from in INAC, the parameter vector  $\mathbf{v}_t$  in RLSEGAC and SEGAC is updated by  $\delta_t^Q$  rather than  $\delta_t^V$ .

**5. Performance Analysis.** In this section, we theoretically analyze the influence of the  $\epsilon$ -greedy Gaussian policy on the performance of compatible AC algorithms.

**5.1. Performance improvement analysis.** In compatible AC algorithms, the traditional Gaussian policy is improved only by policy gradient descent. By contrast, the  $\epsilon$ -greedy Gaussian policy can be indirectly improved by maximizing the advantage functions of candidate actions, even if it does not use policy gradient descent.

**Theorem 5.1.** *For any traditional Gaussian policy  $\vartheta$ , the corresponding  $\epsilon$ -greedy Gaussian policy  $\pi$  with respect to  $Q^\vartheta(\mathbf{s}, \mathbf{a})$  is better than or equal to  $\vartheta$ .*

**Proof:** Let  $Q^\vartheta(\mathbf{s}, \mathbf{a}^*) = \max_{\mathbf{u}_j \in \mathcal{A}(\mathbf{s})} Q^\vartheta(\mathbf{s}, \mathbf{u}_j)$ . Based on the definition and properties of  $\pi$ , for any  $\mathbf{u}_i, \mathbf{u}_j \in \mathcal{A}(\mathbf{s})$ , we have: (i)  $\vartheta(\mathbf{u}_i|\mathbf{s}) = \vartheta(\mathbf{u}_j|\mathbf{s})$ ; (ii) If  $\mathbf{u}_j = \mathbf{a}^*$ , then  $\pi(\mathbf{u}_j|\mathbf{s}) = 2^n (1 - \epsilon + \frac{\epsilon}{2^n}) \vartheta(\mathbf{u}_j|\mathbf{s})$ , else  $\pi(\mathbf{u}_j|\mathbf{s}) = \epsilon \vartheta(\mathbf{u}_j|\mathbf{s})$ . Thus

$$\begin{aligned} \sum_j \pi(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j) &= (1 - \epsilon)2^n \vartheta(\mathbf{a}^*|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{a}^*) + \epsilon \sum_j \vartheta(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j) \\ &\geq (1 - \epsilon) \sum_j \vartheta(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j) + \epsilon \sum_j \vartheta(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j) \\ &= \sum_j \vartheta(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j) \end{aligned}$$

Let  $\pi(\mathbf{s}) \in \mathcal{A}$  denote the action determined by  $\pi$  in  $\mathbf{s}$ . For any  $\mathbf{s} \in \mathcal{S}$ ,

$$\begin{aligned} Q^\vartheta(\mathbf{s}, \pi(\mathbf{s})) &= \int_{\mathcal{A}} \pi(\mathbf{a}|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{a})d\mathbf{a} \\ &= \frac{1}{2^n} \sum_j \int_{\mathcal{A}} \pi(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j)d\mathbf{u}_j \\ &\geq \frac{1}{2^n} \sum_j \int_{\mathcal{A}} \vartheta(\mathbf{u}_j|\mathbf{s})Q^\vartheta(\mathbf{s}, \mathbf{u}_j)d\mathbf{u}_j \\ &= V^\vartheta(\mathbf{s}) \end{aligned}$$

Therefore, by the policy improvement theorem [2], the claim follows.

**5.2. Convergence analysis.** From the definition of the  $\epsilon$ -greedy Gaussian policy, the candidate action set  $\mathcal{A}(\mathbf{s})$  is generated by the traditional Gaussian policy. That means the  $\epsilon$ -greedy Gaussian policy is built on the traditional Gaussian policy, i.e.,  $\pi \sim \vartheta$ . Although its greediness is increased by introducing the  $\epsilon$ -greedy policy, it is still a stochastic policy rather than a deterministic greedy policy. Thus, we believe that the  $\epsilon$ -greedy Gaussian policy does not cause exploration insufficiency and cannot destroy the original convergence of compatible AC algorithms with the traditional Gaussian policy.

**Theorem 5.2.** *For any traditional Gaussian policy  $\vartheta$ , if a compatible AC algorithm with it is convergent, then the algorithm with the corresponding  $\epsilon$ -greedy Gaussian policy  $\pi$  is also convergent.*

**Proof:** Since  $\pi \sim \vartheta$  and  $\nabla_{\boldsymbol{\theta}} \log \pi(\mathbf{a}|\mathbf{s}) = \nabla_{\boldsymbol{\theta}} \log \vartheta(\mathbf{a}|\mathbf{s})$ , the differences between  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\pi$  and  $\nabla_{\boldsymbol{\theta}}^{\text{PG}} J^\vartheta$  are mainly manifest in  $d(\mathbf{s})$  and  $A(\mathbf{s}, \mathbf{a})$ . In other words, the influence of introducing the  $\epsilon$ -greedy policy on the convergence of a compatible AC algorithm is mainly manifest in the critic parameters. By the two-timescale stochastic approximation theory [4], the convergence of the critic parameters is analyzed under a given  $\boldsymbol{\theta}$ . By Theorem 5.1, if  $J^\vartheta$ ,  $\mathbf{v}^\vartheta$  and  $\mathbf{w}^\vartheta$  are convergent under a given  $\boldsymbol{\theta}$ , then  $J^\pi$  is convergent surely, and thus  $\mathbf{v}^\pi$  and  $\mathbf{w}^\pi$  are also convergent surely. The claim follows.

**Corollary 5.1.** *RLSEGAC and SEGAC have a good convergence guarantee.*

**Proof:** It is well known that AC algorithms generally have a reliable convergence guarantee [3, 4], and experimental results in [9, 16] show that LSNAC and SNAC with the traditional Gaussian policy are convergent. Furtherly, RLSEGAC can be viewed as an extension of LSNAC, and SEGAC can also be viewed as an extension of SNAC. By Theorem 5.2, the claim follows.

**6. Experiments.** In this section, in order to demonstrate the effectiveness of the  $\epsilon$ -greedy Gaussian policy and the proposed compatible AC algorithms, we use a continuous mountain-car problem [2, 6] and a continuous puddle-world problem [18, 23] to compare the following algorithm frameworks: (i) RLSEGAC and RLSTAC; (ii) SEGAC and STAC.

**6.1. Experimental settings.** In the mountain-car problem, the task is to learn a throttle control policy, which forces an underpowered car from a valley to the top of the rightmost hill as quickly as possible. At time step  $t$ , the state  $\mathbf{s}_t = (x_t, \dot{x}_t)$  is updated by  $x_{t+1} = \text{bound}[x_t + \dot{x}_{t+1}]$  and  $\dot{x}_{t+1} = \text{bound}[\dot{x}_t + 0.001a_t - 0.0025 \cos(3x_t)]$ , where  $x_t$  and  $\dot{x}_t$  are the position and velocity of the car, the action  $a_t \in [-1, 1]$  is the throttle value, and the bound operation enforces  $x_{t+1} \in [-1.5, 0.5]$  and  $\dot{x}_{t+1} \in [-0.07, 0.07]$ . If  $x_{t+1} = 0.5$ , then the reward  $r_t$  is 100 and the goal is reached; otherwise,  $r_t$  is  $-1$ . In this experiment, a run includes 100 episodes. At each episode, the car starts from an initial state  $\mathbf{s}_0 = (-0.5, 0)$ , and ends when  $t$  reaches 1000 or the car reaches the goal.

In the puddle-world problem, the task is to learn a move policy, which forces an agent from a given position to the top-right corner of the puddle world as quickly as possible. At time step  $t$ , the state  $\mathbf{s}_t = (x_t, y_t) \in [0, 1]^2$  denotes the position in the puddle world, the action  $\mathbf{a}_t = (a_t^x, a_t^y) \in [-0.025, 0.025]^2$  denotes the moves in both dimensions, and the reward for arriving in the successor state  $\mathbf{s}_{t+1} = (p_x, p_y)$  is defined as  $r_t = -1 - 2(\mathcal{N}(p_x, 0.3, 0.1) \cdot \mathcal{N}(p_y, 0.6, 0.03) + \mathcal{N}(p_x, 0.4, 0.03) \cdot \mathcal{N}(p_y, 0.5, 0.1) + \mathcal{N}(p_x, 0.8, 0.03) \cdot \mathcal{N}(p_y, 0.9, 0.1))$ , where  $\mathcal{N}(p_\bullet, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(p_\bullet - \mu)^2}{2\sigma^2}\right)$ . In this experiment, a run includes 100 episodes. At each episode, the agent starts from an initial state  $\mathbf{s}_0 = (0.2, 0.4)$ , and ends when  $t$  reaches 1000 or the goal is reached (i.e.,  $|p_x - 1| + |p_y - 1| < 0.1$ ).

In the implementations of all compatible AC algorithms for both problems, the settings are summarized as follows. (i) The traditional Gaussian policy is defined as  $\vartheta(\mathbf{a}|\mathbf{s}) = \prod_{i=1}^n \mathcal{N}(\boldsymbol{\theta}_i^\top \boldsymbol{\phi}(\mathbf{s}), 1)$ , and the  $\epsilon$ -greedy Gaussian policy  $\pi(\mathbf{a}|\mathbf{s})$  is built on  $\vartheta(\mathbf{a}|\mathbf{s})$  with  $\epsilon = 0.01$ . (ii) The feature vector  $\boldsymbol{\phi}(\mathbf{s})$  is evenly constructed over the state space by  $10 \times 10$  Gaussian RBFs, and the window width of RBFs for the  $i$ -th state component is defined as  $\sigma_i = \frac{ub_i - lb_i}{9\sqrt{2}}$  ( $i = 1, 2$ ), where  $ub_i$  and  $lb_i$  are the upper and lower bounds of the  $i$ -th state component. (iii) All parameter vectors, eligibility traces and  $\bar{r}_0$  are initialized to zero vectors, and the trace-decay parameter  $\lambda$  is set to 0.5. (iv) For all RLSEGAC and RLSTAC algorithms,  $\mathbf{P}_0$  is initialized to  $0.15\mathbf{E}$ . (v) All step sizes are constants, which are selected from  $\{0.00001, 0.00005, 0.0001, \dots, 0.1, 0.5\}$  by using the parameter-sweep method [21]. Similarly, the forgetting factor  $\beta$  is selected from  $\{0.99, 1.0\}$ , and the discount factor  $\gamma$  is selected from  $\{0.97, 0.98, \dots, 1.0\}$  when the compatible AC algorithm uses the discount-return formulation. Considering that  $\beta$  will decay the compatible parameter vector  $\mathbf{w}$  if  $\beta < 1$ , we allow  $\alpha_\epsilon = \alpha_\theta$  during sweeping parameters.

**6.2. Experimental results and analysis.** The best parameters for each compatible AC algorithm and all average performance comparisons over 30 runs are summarized in Tables 1 and 2 and Figures 2 and 3. At the end of each algorithm's name, 'D' or 'A' denotes the discount or average return, and 'R' or 'N' denotes the regular or natural gradient. For example, RLSEGAC-D-R denotes the RLSEGAC algorithm with the discount return and the regular gradient. In Tables 1 and 2, 'Steps-I' and 'Steps-II' denote the average time steps with standard errors to the goal at the 30th episode and the 100th episode. In Figures 2 and 3, steps are the average time steps to the goal as the number of episodes increases.

From Tables 1 and 2 and Figures 2 and 3, we can see the followings. (i) Compared with RLSTAC and STAC algorithms, all RLSEGAC and SEGAC algorithms can obtain better convergence quality. Since all RLSTAC and STAC algorithms also use the best parameters, the performance improvement of RLSEGAC and SEGAC algorithms mainly benefits from the introduction of the  $\epsilon$ -greedy policy rather than the selection of parameters. (ii) The performance improvement of RLSEGAC and SEGAC algorithms is more significant in the puddle-world problem than in the mountain-car problem, since there

TABLE 1. Best parameters and average time steps for the mountain-car problem

|             | $\alpha^{\bar{r}}$ | $\alpha^{\xi}$ | $\alpha^{\theta}$ | $\beta$ | $\gamma$ | Steps-I    | Steps-II   |
|-------------|--------------------|----------------|-------------------|---------|----------|------------|------------|
| RLSTAC-D-R  | –                  | –              | 0.005             | 1.0     | 0.97     | 207.0±49.7 | 126.5±15.3 |
| RLSEGAC-D-R | –                  | –              | 0.01              | 1.0     | 0.98     | 135.9±22.5 | 125.6±26.7 |
| RLSTAC-D-N  | –                  | –              | 0.005             | 1.0     | 0.99     | 133.6±32.6 | 111.7±4.5  |
| RLSEGAC-D-N | –                  | –              | 0.005             | 1.0     | 0.99     | 118.4±13.9 | 111.7±1.7  |
| RLSTAC-A-R  | 0.0001             | –              | 0.01              | 1.0     | –        | 173.9±23.8 | 167.2±22.8 |
| RLSEGAC-A-R | 0.0001             | –              | 0.005             | 1.0     | –        | 167.4±3.2  | 163.7±2.9  |
| RLSTAC-A-N  | 0.0001             | –              | 0.01              | 1.0     | –        | 147.1±33.5 | 129.8±28.9 |
| RLSEGAC-A-N | 0.0001             | –              | 0.01              | 1.0     | –        | 133.0±26.8 | 127.0±19.9 |
| STAC-D-R    | –                  | 0.01           | 0.01              | 0.99    | 1.0      | 195.7±38.0 | 178.8±19.4 |
| SEGAC-D-R   | –                  | 0.01           | 0.01              | 0.99    | 1.0      | 174.9±21.1 | 124.4±19.6 |
| STAC-D-N    | –                  | 0.01           | 0.01              | 0.99    | 0.99     | 191.0±20.3 | 166.4±19.4 |
| SEGAC-D-N   | –                  | 0.01           | 0.01              | 0.99    | 0.99     | 166.4±42.8 | 126.5±29.2 |
| STAC-A-R    | 0.0001             | 0.01           | 0.005             | 1.0     | –        | 200.7±35.0 | 171.6±21.9 |
| SEGAC-A-R   | 0.0001             | 0.01           | 0.01              | 0.99    | –        | 174.6±28.8 | 122.4±17.0 |
| STAC-A-N    | 0.00001            | 0.01           | 0.01              | 0.99    | –        | 182.3±22.3 | 155.5±25.6 |
| SEGAC-A-N   | 0.0001             | 0.01           | 0.005             | 0.99    | –        | 162.4±32.2 | 151.1±37.3 |

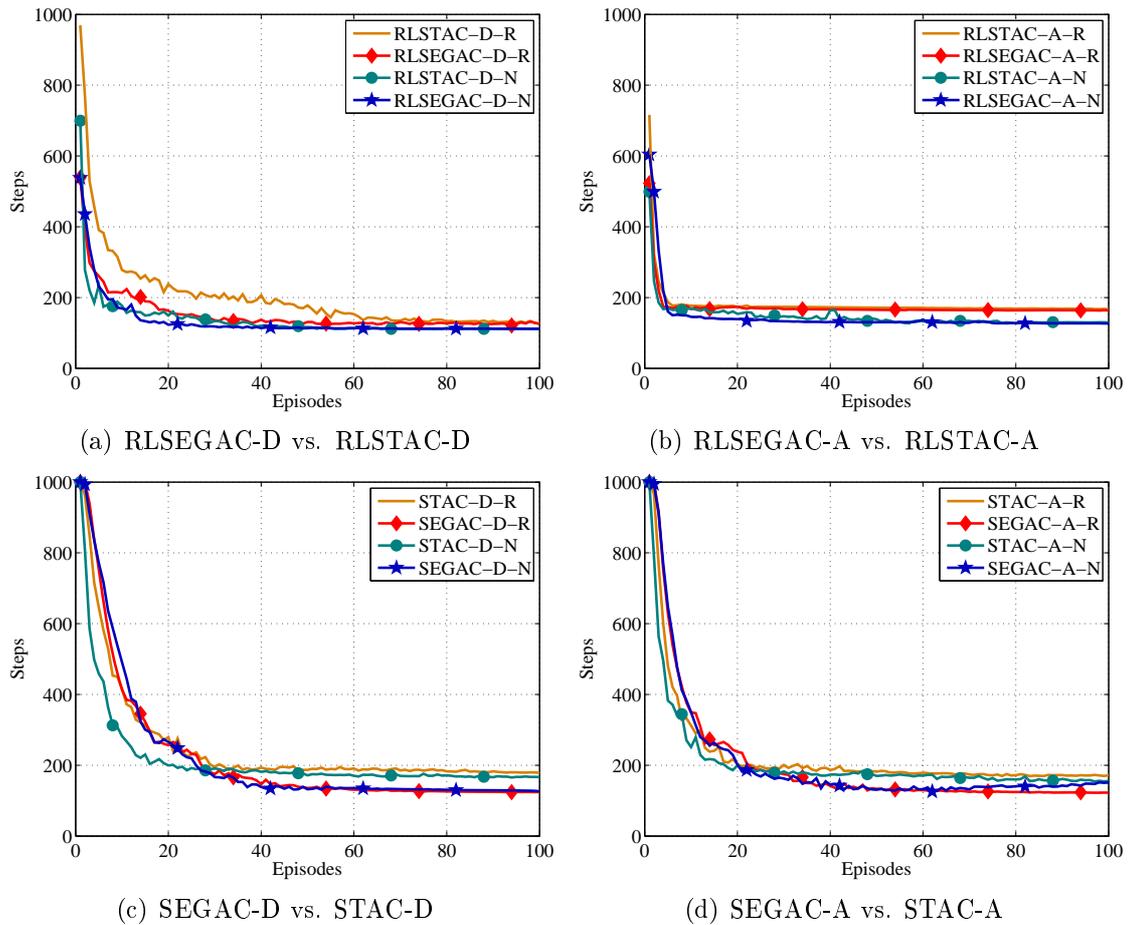


FIGURE 2. Average performance comparisons for the mountain-car problem

TABLE 2. Best parameters and average time steps for the puddle-world problem

|             | $\alpha^{\bar{r}}$ | $\alpha^{\xi}$ | $\alpha^{\theta}$ | $\beta$ | $\gamma$ | Steps-I     | Steps-II  |
|-------------|--------------------|----------------|-------------------|---------|----------|-------------|-----------|
| RLSTAC-D-R  | –                  | –              | 0.005             | 1.0     | 1.0      | 49.6±13.7   | 46.4±8.9  |
| RLSEGAC-D-R | –                  | –              | 0.005             | 1.0     | 1.0      | 46.1±9.7    | 41.5±5.5  |
| RLSTAC-D-N  | –                  | –              | 0.005             | 0.99    | 1.0      | 66.8±10.5   | 61.7±9.5  |
| RLSEGAC-D-N | –                  | –              | 0.005             | 0.99    | 1.0      | 57.3±22.8   | 48.4±6.7  |
| RLSTAC-A-R  | 0.00001            | –              | 0.005             | 1.0     | –        | 118.3±240.7 | 47.0±6.8  |
| RLSEGAC-A-R | 0.00001            | –              | 0.005             | 1.0     | –        | 45.9±9.0    | 41.6±7.3  |
| RLSTAC-A-N  | 0.00001            | –              | 0.005             | 0.99    | –        | 97.2±171.2  | 74.2±61.2 |
| RLSEGAC-A-N | 0.00001            | –              | 0.005             | 0.99    | –        | 56.0±15.6   | 50.6±6.5  |
| STAC-D-R    | –                  | 0.1            | 0.005             | 0.99    | 1.0      | 279.3±404.5 | 45.1±3.3  |
| SEGAC-D-R   | –                  | 0.1            | 0.005             | 0.99    | 1.0      | 106.9±242.9 | 39.8±3.7  |
| STAC-D-N    | –                  | 0.05           | 0.001             | 0.99    | 1.0      | 525.2±453.2 | 47.9±6.6  |
| SEGAC-D-N   | –                  | 0.05           | 0.001             | 0.99    | 1.0      | 69.9±58.6   | 41.4±3.2  |
| STAC-A-R    | 0.00001            | 0.1            | 0.005             | 0.99    | –        | 467.2±460.5 | 51.1±9.5  |
| SEGAC-A-R   | 0.00001            | 0.1            | 0.005             | 0.99    | –        | 137.2±292.5 | 39.1±2.5  |
| STAC-A-N    | 0.00001            | 0.05           | 0.001             | 0.99    | –        | 584.6±453.1 | 48.8±5.2  |
| SEGAC-A-N   | 0.00001            | 0.05           | 0.001             | 0.99    | –        | 67.5±59.9   | 40.4±2.4  |

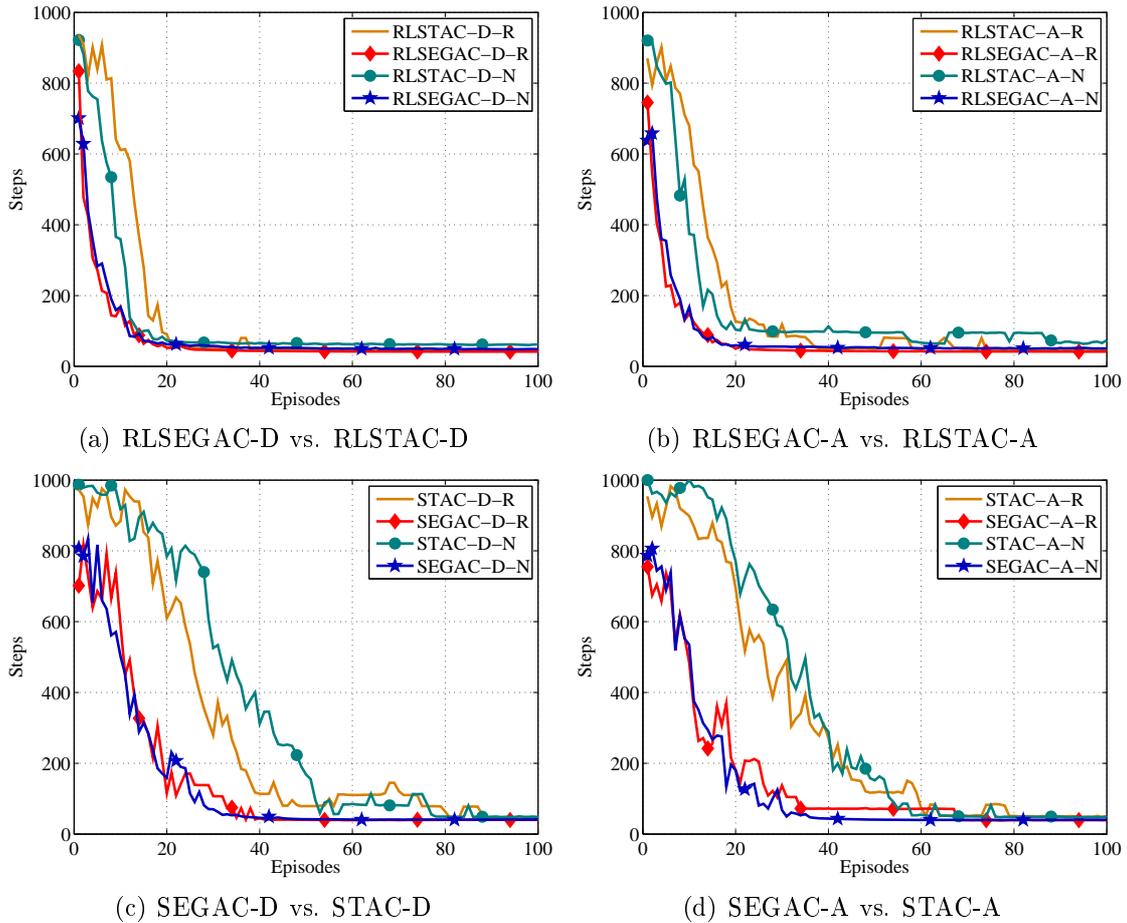


FIGURE 3. Average performance comparisons for the puddle-world problem

are more candidate actions per time step and thus the  $\epsilon$ -greedy Gaussian policy has more greediness. That indicates the  $\epsilon$ -greedy Gaussian policy is more suitable for multidimensional continuous action problems than the traditional Gaussian policy. (iii) In the mountain-car problem, SEGAC algorithms converge slower than STAC algorithms at the first 30 episodes, but they are obviously superior in the final convergence quality. That indicates the  $\epsilon$ -greedy Gaussian policy is helpful to avoid trapping into a local optimal solution. (iv) Compared with SEGAC algorithms, RLSEGAC algorithms converge much faster, since they can reuse samples efficiently [17]. However, SEGAC algorithms are more suitable for real-time RL problems, since they have linear computational complexity per time step.

Compared with the previous work reviewed in Section 1, the  $\epsilon$ -greedy Gaussian policy and the proposed compatible AC frameworks have several other advantages. (i) They are simple and easy to implement. (ii) They do not require extra memory and sampling cost. (iii) The  $\epsilon$ -greedy Gaussian policy bridges the gap between the traditional Gaussian policy and the  $\epsilon$ -greedy policy. It allows us to adjust the greediness of the behaviour policy by selecting  $\epsilon$ . (iv) By using the symmetry of the normal distribution, the  $\epsilon$ -greedy Gaussian policy has the same compatible feature vector as the traditional Gaussian policy, and does not change the update formulas of policy parameters. Thus, it can be used in any compatible AC algorithm supporting continuous states and actions, and can also be used in online, offline, on-policy or off-policy settings. (v) Since the  $\epsilon$ -greedy Gaussian policy is still a stochastic policy and its influence is mainly manifest in policy evaluation, it does not destroy the original convergence of compatible AC algorithms.

**7. Conclusions.** In this paper, we argue that the traditional Gaussian policy has a tendency of over-exploration due to lack of greediness. To overcome this problem, we propose a novel  $\epsilon$ -greedy Gaussian policy, which can be viewed as a hybrid of the traditional Gaussian policy and the  $\epsilon$ -greedy policy. At each time step, it generates a candidate action set by using symmetric Gaussian perturbations, and then uses the  $\epsilon$ -greedy policy to select the behaviour action based on the advantage functions of all candidate actions. For successfully using this policy, we also present RLSEGAC and SEGAC frameworks based on compatible function approximation. Within each framework, we consider the discount return and the average return in the critic, and consider the regular gradient and the natural gradient in the actor. The theoretical analysis shows that the  $\epsilon$ -greedy Gaussian policy is more powerful for policy improvement than the traditional Gaussian policy, and compatible AC algorithms with it do not suffer from divergence although the greediness is increased. Experimental results demonstrate the effectiveness of the proposed policy and frameworks for two benchmark problems with continuous state and action spaces.

This research shows that it is very promising to introduce the  $\epsilon$ -greedy policy into AC algorithms for improving the convergence performance. There are also some interesting topics to be studied in future work. (i) A more thorough experimental study, especially for high-dimensional continuous action problems, should be done. (ii) The comparison of sample efficiency between the  $\epsilon$ -greedy Gaussian policy and the traditional Gaussian policy should be made theoretically. (iii) The  $\epsilon$ -greedy Gaussian policy requires symmetric Gaussian perturbations, which limit the number of candidate actions and also limit the increase of greediness. Thus, other perturbation methods and the corresponding AC algorithms should be further investigated.

**Acknowledgment.** This work is supported in part by the National Natural Science Foundation of China under Grant Nos. 61300192 and 11261015, the Fundamental Research Funds for the Central Universities under Grant No. ZYGX2014J052, and the Natural

Science Foundation of Hainan Province, China, under Grant No. 613153. The authors also gratefully acknowledge the helpful comments and suggestions of the reviewers, which have improved the presentation.

## REFERENCES

- [1] I. Grondman, L. Busoniu, G. A. D. Lopes and R. Babuska, A survey of actor-critic reinforcement learning: Standard and natural policy gradients, *IEEE Trans. Systems, Man, and Cybernetics, Part C: Applications and Reviews*, vol.42, no.6, pp.1291-1307, 2012.
- [2] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, MIT Press, Cambridge, 1998.
- [3] V. R. Konda and J. N. Tsitsiklis, On actor-critic algorithms, *SIAM J. on Control and Optimization*, vol.42, no.4, pp.1143-1166, 2003.
- [4] S. Bhatnagar, R. S. Sutton, M. Ghavamzadeh and M. Lee, Natural actor-critic algorithms, *Technical Report*, Department of Computing Science, University of Alberta, Canada, 2009.
- [5] A. Lazaric, M. Restelli and A. Bonarini, Reinforcement learning in continuous action spaces through sequential Monte Carlo methods, *Proc. of Advances in Neural Information Processing Systems 20*, Vancouver, Canada, pp.833-840, 2007.
- [6] J. A. Martin, J. de Lope and D. Maravall, Robust high performance reinforcement learning through weighted k-nearest neighbors, *Neurocomputing*, vol.74, no.8, pp.1251-1259, 2011.
- [7] J. C. Santamaría, R. S. Sutton and A. Ram, Experiments with reinforcement learning in problems with continuous state and action spaces, *Adaptive Behaviour*, vol.6, no.2, pp.163-217, 1997.
- [8] F. S. Melo and M. Lopes, Fitted natural actor-critic: A new algorithm for continuous state-action MDPs, *Proc. of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Database*, Antwerp, Belgium, pp.66-81, 2008.
- [9] P. Wawrzynski, Real-time reinforcement learning by sequential actor-critics and experience replay, *Neural Networks*, vol.22, no.10, pp.1484-1497, 2009.
- [10] P. Wawrzynski and A. K. Tanwani, Autonomous reinforcement learning with experience replay, *Neural Networks*, vol.41, pp.156-167, 2013.
- [11] H. van Hasselt and M. A. Wiering, Reinforcement learning in continuous action spaces, *Proc. of the IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning*, Honolulu, USA, pp.272-279, 2007.
- [12] H. van Hasselt and M. A. Wiering, Using continuous action spaces to solve discrete problems, *Proc. of the International Joint Conference on Neural Networks*, Atlanta, USA, pp.1149-1156, 2009.
- [13] H. van Hasselt, Reinforcement learning in continuous state and action spaces, *Chap. of Reinforcement Learning: State of the Art*, pp.207-251, 2012.
- [14] N. A. Vien and W. Ertel, Learning via human feedback in continuous state and action spaces, *Proc. of AAAI Fall Symposium on Robots Learning Interactively from Human Teachers*, Arlington, USA, pp.65-72, 2012.
- [15] H. Kimura and S. Kobayashi, An actor-critic algorithm using binary tree action selector, *Trans. of the Society of Instrument and Control Engineers*, vol.4, no.1, pp.1-9, 2007.
- [16] J. Peters and S. Schaal, Natural actor-critic, *Neurocomputing*, vol.71, nos.7-9, pp.1180-1190, 2008.
- [17] Y. Cheng, H. Feng and X. Wang, Efficient data use in incremental actor-critic algorithms, *Neurocomputing*, vol.116, pp.346-354, 2013.
- [18] T. Degris, M. White and R. S. Sutton, Off-policy actor-critic, *Proc. of the 29th International Conference on Machine Learning*, Edinburgh, Scotland, pp.457-464, 2012.
- [19] D. H. Lee and J. J. Lee, Incremental receptive field weighted actor-critic, *IEEE Trans. Industrial Informatics*, vol.9, no.1, pp.62-71, 2013.
- [20] W. Zhu, Y. Jin, Y. Fu and X. Song, Recursive least-squares actor-critic algorithm in continuous space, *Application Research of Computers*, vol.31, no.7, pp.1994-1997, 2014.
- [21] T. Degris, P. M. Pilarski and R. S. Sutton, Model-free reinforcement learning with continuous action in practice, *Proc. of American Control Conference*, Montreal, Canada, pp.2177-2182, 2012.
- [22] T. Rückstieß, M. Felder and J. Schmidhuber, State-dependent exploration for policy gradient methods, *Proc. of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Database*, Antwerp, Belgium, pp.234-249, 2008.
- [23] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra and M. Riedmiller, Deterministic policy gradient algorithms, *Proc. of the 31st International Conference on Machine Learning*, Beijing, China, pp.387-395, 2014.

- [24] C. Zhang and Q. Zhu, Actor-critic algorithms based on symmetric perturbation sampling, *Control and Decision*, vol.30, no.12, pp.2161-2167, 2015.
- [25] R. S. Sutton, D. McAllester, S. Singh and Y. Mansour, Policy gradient methods for reinforcement learning with function approximation, *Proc. of Advances in Neural Information Processing Systems 12*, Denver, USA, pp.1057-1063, 1999.
- [26] M. Geist and O. Pietquin, Revisiting natural actor-critics with value function approximation, *Proc. of the 7th International Conference on Modeling Decisions for Artificial Intelligence*, Perpignan, France, pp.207-218, 2010.
- [27] J. N. Tsitsiklis and B. V. Roy, An analysis of temporal-difference learning with function approximation, *IEEE Trans. Automatic Control*, vol.42, no.5, pp.674-690, 1997.
- [28] J. N. Tsitsiklis and B. V. Roy, Average cost temporal-difference learning, *Automatica*, vol.35, no.11, pp.1799-1808, 1999.
- [29] N. Heess, D. Silver and Y. W. Teh, Actor-critic reinforcement learning with energy-based policies, *Proc. of the 10th European Workshop on Reinforcement Learning*, Edinburgh, Scotland, pp.43-58, 2012.
- [30] X. Xu, H. G. He and D. Hu, Efficient reinforcement learning using recursive least-squares methods, *Journal of Artificial Intelligence Research*, vol.16, pp.259-292, 2002.
- [31] T. Mori, Y. Nakamura and S. Ishii, Off-policy natural actor-critic, *Technical Report 2005007*, Nara Institute of Science and Technology, Japan, 2005.
- [32] J. Park, J. Kim and D. Kang, An RLS-based natural actor-critic algorithm for locomotion of a two-linked robot arm, *Proc. of International Conf. Computational Intelligence and Security*, Xi'an, China, pp.65-72, 2005.
- [33] T. Morimura, E. Uchibe and K. Doya, Utilizing natural gradient in temporal difference reinforcement learning with eligibility traces, *Proc. of the 2nd International Symposium on Information Geometry and its Applications*, Tokyo, Japan, pp.256-263, 2005.