
Discrete Action On-Policy Learning with Action-Value Critic

Yuguang Yue
UT-Austin

Yunhao Tang
Columbia University

Mingzhang Yin
UT-Austin

Mingyuan Zhou
UT-Austin

Abstract

Reinforcement learning (RL) in discrete action space is ubiquitous in real-world applications, but its complexity grows exponentially with the action-space dimension, making it challenging to apply existing on-policy gradient based deep RL algorithms efficiently. To effectively operate in multidimensional discrete action spaces, we construct a critic to estimate action-value functions, apply it on correlated actions, and combine these critic estimated action values to control the variance of gradient estimation. We follow rigorous statistical analysis to design how to generate and combine these correlated actions, and how to sparsify the gradients by shutting down the contributions from certain dimensions. These efforts result in a new discrete action on-policy RL algorithm that empirically outperforms related on-policy algorithms relying on variance control techniques. We demonstrate these properties on OpenAI Gym benchmark tasks, and illustrate how discretizing the action space could benefit the exploration phase and hence facilitate convergence to a better local optimal solution thanks to the flexibility of discrete policy.

1 INTRODUCTION

There has been significant recent interest in using model-free reinforcement learning (RL) to address complex real-world sequential decision making tasks [Silver et al., 2018, MacAlpine and Stone, 2017, OpenAI, 2018]. With the help of deep neural networks, model-free deep RL algorithms have been successfully implemented in a variety of tasks, including game playing [Silver et al.,

2016, Mnih et al., 2013] and robotic controls [Levine et al., 2016]. Among those model-free RL algorithms, policy gradient (PG) algorithms are a class of methods that parameterize the policy function and apply gradient-based methods to make updates. It has been shown to succeed in solving a range of challenging RL tasks [Mnih et al., 2016, Schulman et al., 2015a, Lillicrap et al., 2015, Schulman et al., 2017, Wang et al., 2016, Haarnoja et al., 2018, Liu et al., 2017b]. Despite directly targeting at maximizing the expected rewards, PG suffers from problems including having low sample efficiency [Haarnoja et al., 2018] for on-policy PG algorithms and undesirable sensitivity to hyper-parameters for off-policy algorithms [Lillicrap et al., 2015].

On-policy RL algorithms use on-policy samples to estimate the gradients for policy parameters, as routinely approximated by Monte Carlo (MC) estimation that often comes with large variance. A number of techniques have sought to alleviate this problem for continuous action spaces [Gu et al., 2016, Grathwohl et al., 2017, Liu et al., 2017a, Wu et al., 2018], while relatively fewer have been proposed for discrete action spaces [Grathwohl et al., 2017, Yin et al., 2019]. In practice, RL with discrete action space is ubiquitous in fields including recommendation system [Dulac-Arnold et al., 2015], bidding system [Hu et al., 2018], gaming [Mnih et al., 2013], to name a few. It plays an important role in the early stage of RL development [Sutton and Barto, 2018], and many value-based algorithms [Watkins and Dayan, 1992, Mnih et al., 2013, Van Hasselt et al., 2016] can handle such setup when the action space is not large. However, when the action space is multidimensional, the number of unique actions grows exponentially with the dimension, leading to an intractable combinatorial optimization problem at every single step that prevents the application of most value-based RL methods.

Under the setting of high-dimensional discrete action space, policy-gradient based algorithms can still be applied if we assume the joint distribution over discrete actions to be factorized across dimensions, so that the joint policy is still tractable [Jaśkowski et al., 2018, Andrychowicz et al., 2018]. Then the challenge boils down to obtaining a gradient estimator that can well

control its variance. Though many variance reduction techniques have been proposed for discrete variables [Jang et al., 2016, Tucker et al., 2017, Yin and Zhou, 2019, Raiko et al., 2014], they either provide biased gradients or are not applicable to multidimensional RL settings.

In this paper, we propose Critic-ARSM (CARSM) policy gradient, which improves the recently proposed augment-REINFORCE-swap-merge (ARSM) gradient estimator of Yin et al. [2019] and integrates it with action-value function evaluation, to accomplish three-fold effects: **1)** CARSM sparsifies the ARSM gradient and introduces an action-value Critic to work with multidimensional discrete actions spaces; **2)** By estimating the rewards of a set of correlated discrete actions via the proposed action-value Critic, and combining these rewards for variance reduction, CARSM achieves better sample efficiency compared with other variance-control methods such as A2C [Mnih et al., 2016] and RELAX [Grathwohl et al., 2017]; **3)** CARSM can be easily applied to other RL algorithms using REINFORCE or its variate as the gradient estimator. Although we mainly focus on on-policy algorithms, our algorithm can also be potentially applied to off-policy algorithms with the same principle; we leave this extension for future study.

The paper proceeds as follows. In Section 2, we briefly review existing on-policy learning frameworks and variance reduction techniques for discrete action space. In Section 3, we introduce CARSM from both theoretical and practical perspectives. In Section 4, we first demonstrate the potential benefits of discretizing a continuous control task compared with using a diagonal Gaussian policy, then show the high sample efficiency of CARSM from an extensive range of experiments and illustrate that CARSM can be plugged into state-of-arts on-policy RL learning frameworks such as Trust Region Policy Optimization (TRPO) [Schulman et al., 2015a]. Python (TensorFlow) code is available at <https://github.com/yuguangyue/CARSM>.

2 PRELIMINARIES

RL is often formulated as learning under a Markov decision process (MDP). Its action space \mathcal{A} is dichotomized into either discrete (*e.g.*, $\mathcal{A} = \{1, \dots, 100\}$) or continuous (*e.g.*, $\mathcal{A} = [-1, 1]$). In an MDP, at discrete time $t \geq 0$, an agent in state $\mathbf{s}_t \in \mathcal{S}$ takes action $a_t \in \mathcal{A}$, receives instant reward $r(\mathbf{s}_t, a_t) \in \mathbb{R}$, and transits to next state $\mathbf{s}_{t+1} \sim \mathcal{P}(\cdot | \mathbf{s}_t, a_t)$. Let $\pi : \mathcal{S} \mapsto \mathcal{P}(\mathcal{A})$ be a mapping from the state to a distribution over actions. We define the expected cumulative rewards under π as

$$J(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, a_t) \right], \quad (1)$$

where $\gamma \in (0, 1]$ is a discount factor. The objective of RL is to find the (sub-)optimal policy $\pi^* = \arg \max_\pi J(\pi)$. In practice, it is infeasible to search through all policies and hence one typically resorts to parameterizing the policy π_θ with θ .

2.1 On-Policy Optimization

We introduce on-policy optimization methods from a constrained optimization point of view to unify the algorithms we will discuss in this article. In practice, we want to solve the following constrained optimization problem as illustrated in Schulman et al. [2015a]:

$$\begin{aligned} \max_{\theta} \quad & \mathbb{E}_{\pi_{\theta_{\text{old}}}} \left[\frac{\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_{\theta_{\text{old}}}(\mathbf{a}_t | \mathbf{s}_t)} Q^{\pi_{\theta_{\text{old}}}}(\mathbf{s}_t, \mathbf{a}_t) \right] \\ \text{subject to} \quad & D(\theta_{\text{old}}, \theta) \leq \epsilon, \end{aligned}$$

where $Q^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) = \mathbb{E}_{\pi_\theta} [\sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}'_t)]$ is the action-value function and $D(\cdot, \cdot)$ is some metric that measures the closeness between θ_{old} and θ .

A2C Algorithm: One choice of $D(\cdot, \cdot)$ is the L_2 norm, which will lead us to first-order gradient ascent. By applying first-order Taylor expansion on π_θ around θ_{old} , the problem can be re-written as maximizing $\mathbb{E}_{\pi_{\theta_{\text{old}}}} [Q^{\pi_{\theta_{\text{old}}}}(\mathbf{s}_t, \mathbf{a}_t)] + \nabla_{\theta} J(\pi_{\theta_{\text{old}}})^T (\theta - \theta_{\text{old}})$ subject to $\|\theta - \theta_{\text{old}}\|_2 \leq \epsilon$, which will result in a gradient ascent update scheme; note $\nabla_{\theta} J(\pi_{\theta_{\text{old}}}) := \nabla_{\theta} J(\pi_{\theta})|_{\theta=\theta_{\text{old}}}$. Based on REINFORCE [Williams, 1992], the gradient of the original objective function (1) can be written as

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right]. \quad (2)$$

However, a naive Monte Carlo estimation of (2) has large variance that needs to be controlled. A2C algorithm [Mnih et al., 2016] adds value function $V^{\pi_{\theta}}(\mathbf{s}) := \mathbb{E}_{a_t \sim \pi_{\theta}} [Q^{\pi_{\theta}}(\mathbf{s}_t, a_t)]$ as a baseline and obtains a low-variance estimator of $\nabla_{\theta} J(\pi_{\theta})$ as

$$\mathbf{g}_{\text{A2C}} = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^{\infty} A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right], \quad (3)$$

where $A^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) = Q^{\pi_{\theta}}(\mathbf{s}_t, \mathbf{a}_t) - V^{\pi_{\theta}}(\mathbf{s}_t)$ is called the Advantage function.

Trust Region Policy Optimization: The other choice of metric $D(\cdot, \cdot)$ could be KL-divergence, and the update from this framework is introduced as TRPO [Schulman et al., 2015a]. In practice, this constrained optimization problem is reformulated as follows:

$$\begin{aligned} \max_{\theta} \quad & \nabla_{\theta} J(\pi_{\theta_{\text{old}}})^T (\theta - \theta_{\text{old}}) \\ \text{subject to} \quad & \frac{1}{2} (\theta_{\text{old}} - \theta)^T H (\theta_{\text{old}} - \theta) \leq \delta, \end{aligned}$$

where H is the second-order derivative $\nabla_{\theta}^2 D_{\text{KL}}(\theta_{\text{old}} | | \theta)|_{\theta=\theta_{\text{old}}}$. An analytic update step for this optimization problem can be expressed as

$$\theta = \theta_{\text{old}} + \sqrt{\frac{2\delta}{\mathbf{d}^T H^{-1} \mathbf{d}}} \mathbf{d}, \quad (4)$$

where $\mathbf{d} = H^{-1} \nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}_{\text{old}}})$, and in practice the default choice of $\nabla_{\boldsymbol{\theta}} J(\pi_{\boldsymbol{\theta}_{\text{old}}})$ is \mathbf{g}_{A2C} as defined at (3).

2.2 Variance Control Techniques

Besides the technique of using state-dependent baseline to reduce variance as in (3), two recent works propose alternative methods for variance reduction in discrete action space settings. For the sake of space, we defer to Grathwohl et al. [2017] for the detail about the RELAX algorithm and briefly introduce ARSM here.

ARSM Policy Gradient: The ARSM gradient estimator can be used to backpropagate unbiased and low-variance gradients through a sequence of unidimensional categorical variables [Yin et al., 2019]. It comes up with a reparametrization formula for discrete random variable, and combines it with a parameter-free self-adjusted baseline to achieve variance reduction.

Instead of manipulating on policy parameters $\boldsymbol{\theta}$ directly, ARSM turns to reduce variance on the gradient with respect to the logits $\boldsymbol{\phi}$, before backpropagating it to $\boldsymbol{\theta}$ using the chain rule. Let us assume

$$\pi_{\boldsymbol{\theta}}(a_t | \mathbf{s}_t) = \text{Categorical}(a_t | \sigma(\boldsymbol{\phi}_t)), \quad \boldsymbol{\phi}_t := \mathcal{T}_{\boldsymbol{\theta}}(\mathbf{s}_t),$$

where $\sigma(\cdot)$ denotes the softmax function and $\mathcal{T}_{\boldsymbol{\theta}}(\cdot) \in \mathbb{R}^C$ denotes a neural network, which is parameterized by $\boldsymbol{\theta}$ and has an output dimension of C .

Denote $\boldsymbol{\varpi}^{c \Leftarrow j}$ as the vector obtained by swapping the c th and j th elements of vector $\boldsymbol{\varpi}$, which means $\varpi_j^{c \Leftarrow j} = \varpi_c$, $\varpi_c^{c \Leftarrow j} = \varpi_j$, and $\varpi_i^{c \Leftarrow j} = \varpi_i$ if $i \notin \{c, j\}$. Following the derivation from Yin et al. [2019], the gradient with respect to ϕ_{tc} can be expressed as

$$\begin{aligned} \nabla_{\phi_{tc}} J(\boldsymbol{\phi}_{0:\infty}) &= \mathbb{E}_{\mathcal{P}(\mathbf{s}_t | \mathbf{s}_0, \pi_{\boldsymbol{\theta}}) \mathcal{P}(\mathbf{s}_0)} \left\{ \gamma^t \mathbb{E}_{\boldsymbol{\varpi}_t \sim \text{Dir}(\mathbf{1}_C)} [g_{tc}] \right\}, \\ g_{tc} &:= \sum_{j=1}^C \left[Q(\mathbf{s}_t, a_t^{c \Leftarrow j}) - \frac{1}{C} \sum_{m=1}^C Q(\mathbf{s}_t, a_t^{m \Leftarrow j}) \right] \left(\frac{1}{C} - \varpi_{tj} \right), \end{aligned}$$

where $\mathcal{P}(\mathbf{s}_t | \mathbf{s}_0, \pi_{\boldsymbol{\theta}})$ is the marginal form of $\prod_{t'=0}^{t-1} \mathcal{P}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, a_{t'}) \text{Categorical}(a_{t'}; \sigma(\boldsymbol{\phi}_{t'}))$, $\boldsymbol{\varpi}_t \sim \text{Dir}(\mathbf{1}_C)$, and $a_t^{c \Leftarrow j} := \arg \min_{i \in \{1, \dots, C\}} \varpi_{ti}^{c \Leftarrow j} e^{-\phi_{ti}}$. In addition, $a_t^{c \Leftarrow j}$ is called a *pseudo action* to differentiate it from the true action $a_t := \arg \min_{i \in \{1, \dots, C\}} \varpi_{ti} e^{-\phi_{ti}}$.

Applying the chain rule leads to ARSM policy gradient:

$$\begin{aligned} \mathbf{g}_{\text{ARSM}} &= \sum_{t=0}^{\infty} \sum_{c=1}^C \frac{\partial J(\boldsymbol{\phi}_{0:\infty})}{\partial \phi_{tc}} \frac{\partial \phi_{tc}}{\partial \boldsymbol{\theta}} \\ &= \mathbb{E}_{\mathbf{s}_t \sim \rho_{\pi, \gamma}(\mathbf{s})} \left\{ \mathbb{E}_{\boldsymbol{\varpi}_t \sim \text{Dir}(\mathbf{1}_C)} \left[\nabla_{\boldsymbol{\theta}} \sum_{c=1}^C g_{tc} \phi_{tc} \right] \right\}, \end{aligned}$$

where $\rho_{\pi, \gamma}(\mathbf{s}) := \sum_{t=0}^{\infty} \gamma^t \mathcal{P}(\mathbf{s}_t = \mathbf{s} | \mathbf{s}_0, \pi_{\boldsymbol{\theta}})$ is the unnormalized discounted state visitation frequency.

In Yin et al. [2019], $Q(\mathbf{s}_t, a_t^{c \Leftarrow j})$ are estimated by MC integration, which requires multiple MC rollouts at each timestep if there are *pseudo actions* that differ from

the true action. This estimation largely limits the implementation of ARSM policy gradient to small action space due to the high computation cost. The maximal number of unique *pseudo actions* grows quadratically with the number of actions along each dimension and a long episodic task will result in more MC rollouts too. To differentiate it from the new algorithm, we refer to it as ARSM-MC.

3 CARSM POLICY GRADIENT

In this section, we introduce Critic-ARSM (CARSM) policy gradient for multidimensional discrete action space. CARSM improves ARSM-MC in the following two aspects: 1. ARSM-MC only works for unidimensional RL settings while CARSM generalizes it to multidimensional ones with sparsified gradients. 2. CARSM can be applied to more complicated tasks as it employs an action-value function critic to remove the need of running multiple MC rollouts for a single estimation, which largely improves the sample efficiency.

For an RL task with K -dimensional C -way discrete action space, we assume different dimensions $a_{tk} \in \{1, \dots, C\}$ of the multidimensional discrete action $\mathbf{a}_t = (a_{t1}, \dots, a_{tK})$ are independent given logits $\boldsymbol{\phi}_t$ at time t , that is $a_{t1} \perp a_{t2} \cdots \perp a_{tK} | \boldsymbol{\phi}_t$. For the logit vector $\boldsymbol{\phi}_t \in \mathbb{R}^{KC}$, which can be decomposed as $\boldsymbol{\phi}_t = (\boldsymbol{\phi}'_{t1}, \dots, \boldsymbol{\phi}'_{tK})'$, $\boldsymbol{\phi}_{tk} = (\phi_{tk1}, \dots, \phi_{tkC})'$, we assume

$$P(\mathbf{a}_t | \boldsymbol{\phi}_t) = \prod_{k=1}^K \text{Categorical}(a_{tk}; \sigma(\boldsymbol{\phi}_{tk})).$$

Theorem 1 (Sparse ARSM for multidimensional discrete action space). *The element-wise gradient of $J(\boldsymbol{\phi}_{0:\infty})$ with respect to ϕ_{tkc} can be expressed as*

$$\nabla_{\phi_{tkc}} J(\boldsymbol{\phi}_{0:\infty}) = \mathbb{E}_{\mathcal{P}(\mathbf{s}_t | \mathbf{s}_0, \pi_{\boldsymbol{\theta}}) \mathcal{P}(\mathbf{s}_0)} \left\{ \gamma^t \mathbb{E}_{\boldsymbol{\varpi}_t \sim \prod_{k=1}^K \text{Dir}(\boldsymbol{\varpi}_{tk}; \mathbf{1}_C)} [g_{tkc}] \right\},$$

where $\boldsymbol{\varpi}_{tk} = (\varpi_{tk1}, \dots, \varpi_{tkC})' \sim \text{Dir}(\mathbf{1}_C)$ is the Dirichlet random vector for dimension k , state t and

$$g_{tkc} := \begin{cases} 0, & \text{if } a_{tk}^{c \Leftarrow j} = a_{tk} \text{ for all } (c, j) \\ \sum_{j=1}^C [\Delta_{c,j}(\mathbf{s}_t, \mathbf{a}_t)] \left(\frac{1}{C} - \varpi_{tkj} \right), & \text{otherwise} \end{cases}$$

$$\Delta_{c,j}(\mathbf{s}_t, \mathbf{a}_t) := Q(\mathbf{s}_t, \mathbf{a}_t^{c \Leftarrow j}) - \frac{1}{C} \sum_{m=1}^C Q(\mathbf{s}_t, \mathbf{a}_t^{m \Leftarrow j}),$$

$$\mathbf{a}_t^{c \Leftarrow j} := (a_{t1}^{c \Leftarrow j}, \dots, a_{tK}^{c \Leftarrow j})',$$

$$a_{tk}^{c \Leftarrow j} := \arg \min_{i \in \{1, \dots, C\}} \varpi_{tki}^{c \Leftarrow j} e^{-\phi_{tki}}.$$

We defer the proof to Appendix A. One difference from the original ARSM [Yin et al., 2019] is the values of g_{tkc} , where we obtain a sparse estimation that shuts down the k th dimension if $a_{tk}^{c \Leftarrow j} = a_{tk}$ for all (c, j) and hence there is no more need to calculate $\nabla_{\boldsymbol{\theta}} \phi_{tkc}$ for all c belonging to dimension k at time t . One immediate benefit from this sparse gradient estimation is to reduce the noise from that specific dimension because the Q

function is always estimated with either MC estimation or Temporal Difference (TD) [Sutton and Barto, 2018], which will introduce variance and bias, respectively.

In ARSM-MC, the action-value function is estimated by MC rollouts. Though it returns unbiased estimation, it inevitably decreases the sample efficiency and prevents it from applying to more sophisticated tasks. Therefore, CARSM proposes using an action-value function critic \hat{Q}_ω parameterized by ω to estimate the Q function. Replacing Q with \hat{Q}_ω in Theorem 1, we obtain $\hat{g}_{tkc} \approx g_{tkc}$ as the empirical estimation of $\nabla_{\phi_{tkc}} J(\pi)$, and hence the CARSM estimation for $\nabla_{\theta} J(\pi) = \sum_{t=0}^{\infty} \sum_{k=1}^K \sum_{c=1}^C \frac{\partial J(\phi_{0:\infty})}{\phi_{tkc}} \frac{\phi_{tkc}}{\partial \theta}$ becomes

$$\hat{g}_{\text{CARSM}} = \nabla_{\theta} \sum_t \sum_{k=1}^K \sum_{c=1}^C \hat{g}_{tkc} \phi_{tkc}.$$

Note the number of unique values in $\{\mathbf{a}_t^{c=j}\}_{c,j}$ that differ from the true action \mathbf{a}_t is always between 0 and $C(C-1)/2 - 1$, regardless of how large K is. The dimension shutdown property further sparsifies the gradients, removing the noise of the dimensions that have no *pseudo actions*.

Design of Critic: A practical challenge of CARSM is that it is notoriously hard to estimate action-value functions for on-policy algorithm because the number of samples are limited and the complexity of the action-value function quickly increases with dimension K . A natural way to overcome the limitation of samples is the reuse of historical data, which has been successfully implemented in previous studies [Gu et al., 2016, Lillicrap et al., 2015]. The idea is to use the transitions $\{\mathbf{s}_\ell, r_\ell, \mathbf{a}_\ell, \mathbf{s}'_\ell\}$'s from the replay buffer to construct target values for the action-value estimator under the current policy. More specifically, we can use one-step TD to rewrite the target value of critic \hat{Q}_ω network with these off-policy samples as

$$y_\ell^{\text{off}} = r(\mathbf{s}_\ell, \mathbf{a}_\ell) + \gamma \mathbb{E}_{\tilde{\mathbf{a}} \sim \pi(\cdot | \mathbf{s}'_\ell)} \hat{Q}_\omega(\mathbf{s}'_\ell, \tilde{\mathbf{a}}), \quad (5)$$

where the expectation part can be evaluated with either an exact computation when the action space size C^K is not large, or with MC integration by drawing random samples from $\tilde{\mathbf{a}} \sim \pi_\theta(\cdot | \mathbf{s}'_\ell)$ and averaging $\hat{Q}_\omega(\mathbf{s}'_\ell, \tilde{\mathbf{a}})$ over these random samples. This target value only uses one-step estimation, and can be extended to n -step TD by adding additional importance sampling weights.

Since we have on-policy samples, it is natural to also include them to construct unbiased targets for $\hat{Q}_\omega(\mathbf{s}_t, \mathbf{a}_t)$:

$$y_t^{\text{on}} = \sum_{t'=t}^{\infty} \gamma^{t'-t} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}).$$

Then we optimize parameters ω by minimizing the Bellman error between the targets and critic as

$$\sum_{\ell=0}^L [y_\ell^{\text{off}} - \hat{Q}_\omega(\mathbf{s}_\ell, \mathbf{a}_\ell)]^2 + \sum_{t=0}^T [y_t^{\text{on}} - \hat{Q}_\omega(\mathbf{s}_t, \mathbf{a}_t)]^2,$$

where L is the number of off-policy samples and T is the number of on-policy samples. In practice, the performance varies with the ratio between L and T , which reflects the trade-off between bias and variance. We choose $L = T$, which is found to achieve good performance across all tested RL tasks.

Target network update: Another potential problem of CARSM is the dependency between the action-value function and policy. Though CARSM is a policy-gradient based algorithm, the gradient estimation procedure is closely related with the action-value function, which may lead to divergence of the estimation as mentioned in previous studies [Mnih et al., 2016, Lillicrap et al., 2015, Bhatnagar et al., 2009, Maei et al., 2010]. Fortunately, this issue has been addressed, to some extent, with the help of target network update [Mnih et al., 2013, Lillicrap et al., 2015], and we borrow that idea into CARSM for computing policy gradient. In detail, we construct two target networks corresponding to the policy network and Q critic network, respectively; when computing the target of critic network in (5), instead of using the current policy network and Q critic, we use a smoothed version of them to obtain the target value, which can be expressed as

$$y_\ell^{\text{off}} = r(\mathbf{s}_\ell, \mathbf{a}_\ell) + \gamma \mathbb{E}_{\tilde{\mathbf{a}} \sim \pi'(\cdot | \mathbf{s}'_\ell)} Q'_\omega(\mathbf{s}'_\ell, \tilde{\mathbf{a}}),$$

where π' and Q'_ω denote the target networks. These target networks are updated every episode in a ‘‘soft’’ update manner, as in Lillicrap et al. [2015], by

$$\omega^{Q'} \leftarrow \tau \omega^Q + (1 - \tau) \omega^{Q'}, \quad \theta^{\pi'} \leftarrow \tau \theta^\pi + (1 - \tau) \theta^{\pi'},$$

which is an exponential moving average of the policy network and action value function network parameters, with τ as the smoothing parameter.

Annealing on entropy term: In practice, maximizing the maximum entropy (ME) objective with an annealing coefficient is often a good choice to encourage exploration and achieve a better sub-optimal solution, and the CARSM gradient estimator for ME would be

$$\mathbf{g}_{\text{CARSM}}^{\text{ME}} = \mathbf{g}_{\text{CARSM}} + \lambda \sum_{k=1}^K \nabla_{\theta} \mathbb{H}(\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)),$$

where $\mathbb{H}(\cdot)$ denotes the entropy term and λ is the annealing coefficient. The entropy term can be expressed explicitly because π is factorized over its dimensions and there are finite actions along each dimension.

Delayed update: As an accurate critic plays an important role for ARSM to estimate gradient, it would be helpful to adopt the delayed update trick of Fujimoto et al. [2018]. In practice, we update the critic network several times before updating the policy network.

In addition to the Python (TensorFlow) code in the Supplementary Material, we also provide detailed pseudo code to help understand the implementation of CARSM in Appendix C.

4 EXPERIMENTS

Our experiments aim to answer the following questions: **(a)** How does the proposed CARSM algorithm perform when compared with ARSM-MC (when ARSM-MC is not too expensive to run)? **(b)** Is CARSM able to efficiently solve tasks with a large discrete action space? **(c)** Does CARSM have better sample efficiency than the algorithms, such as A2C and RELAX, that have the same idea of using baselines for variance reduction? **(d)** Can CARSM be integrated into more sophisticated RL learning frameworks such as TRPO to achieve an improved performance? Since we run trials on some discretized continuous control tasks, another fair question would be: **(e)** Will discretization help learning? If so, what are possible explanations?

We consider benchmark tasks provided by OpenAI Gym classic-control and MuJoCo simulators [Todorov et al., 2012]. We compare the proposed CARSM with ARSM-MC [Yin et al., 2019], A2C [Mnih et al., 2016], and RELAX [Grathwohl et al., 2017]; all of them rely on introducing baseline functions to reduce gradient variance, making it fair to compare them against each other. We then integrate CARSM into TRPO by replacing its A2C gradient estimator for $\nabla_{\theta} J(\theta)$. Performance evaluation show that a simple plug-in of CARSM estimator can bring the improvement. Details on experimental settings can be found in Appendix B.2.

On our experiments with tasks in continuous control domain, we discretize the continuous action space uniformly to get a discrete action space. More specifically, if the action space is $\mathcal{A} = [-1, 1]^K$, and we discretize it to C actions at each dimension, the action space would become $\tilde{\mathcal{A}} = \left\{ \frac{-C+1}{C-1}, \frac{-C+3}{C-1}, \dots, \frac{C-1}{C-1} \right\}^K$.

There are two motivations of discretizing the action space. First, MuJoCo tasks are a set of standard comparable tasks that naturally have multidimensional action spaces, which is the case we are interested in for CARSM. Second, as illustrated in Tang and Agrawal [2019], discrete policy is often more expressive than diagonal-Gaussian policy, leading to better exploration. We will illustrate this point by experiments.

4.1 Motivation and Illustration

One distinction between discrete and Gaussian policies is that a discrete policy can learn multi-modal and skewed distributions while a Gaussian policy can only support uni-modal, symmetric, and bell-shaped distributions. This intrinsic difference could lead to significantly difference on exploration, as reflected by the toy example presented below, which will often lead to different sub-optimal solutions in practice.

To help better understand the connections between

multi-modal policy and exploration, we take a brief review of RL objective function from an energy-based distribution point of view. For a bandit problem with reward function $r(a) : \mathcal{A} \rightarrow \mathbb{R}$, we denote the true reward induced distribution as $p(a) \propto e^{r(a)}$. The objective function in (1) can be reformulated as

$$\mathbb{E}_{a \sim \pi_{\theta}(a)}[r(a)] = -\text{KL}(\pi_{\theta}(a) || p(a)) - \mathbb{H}(\pi_{\theta}).$$

The KL-divergence term matches the objective function of variational inference (VI) [Blei et al., 2017] in approximating $p(a)$ with distribution $\pi_{\theta}(a)$, while the second term is the entropy of policy π_{θ} . Therefore, if we use maximum entropy objective [Haarnoja et al., 2017], which is maximizing $\mathbb{E}_{a \sim \pi_{\theta}(a)}[r(a)] + \mathbb{H}(\pi_{\theta})$, we will get an VI approximate solution. Suppose $p(a)$ is a multi-modal distribution, due to the inherent property of VI [Blei et al., 2017], if π_{θ} is a Gaussian distribution, it will often underestimate the variance of $p(a)$ and capture only one density mode. By contrast, if π_{θ} is a discrete distribution, it can capture the multi-modal property of $p(a)$, which will lead to more exploration before converging to a more deterministic policy.

We design a simple toy example to reflect these differences. We restrict the action space to $[-1, 1]$, and the true reward function is a concatenation of two quadratic functions (as shown in Figure 1 left panel red curves) that intersect at a middle point m . We fix the left sub-optimal point as the global optimal one and control the position of m to get tasks with various difficulty levels. More specifically, the closer m to -1 , the more explorations needed to converge to the global optimal. We defer the detailed experiment setting to Appendix B.1. We run 100 trials of both Gaussian policy and discrete policy and show their behaviors.

We show, in Figure 1 left panel, the learning process of both Gaussian and discrete policies with a quadratic annealing coefficient for the entropy term, and, in right panel, a heatmap where each entry indicates the average density of each action at one iteration. In this case where $m = -0.8$, the signal from the global optimal point has a limited range which requires more explorations during the training process. Gaussian policy can only explore with unimodal distribution and fail to capture the global optimal all the time. By contrast, discrete policy can learn the bi-modal distribution in the early stage, gradually concentrate on both the optimal and sub-optimal peaks before collecting enough samples, and eventually converges to the optimal peak. More explanations can be found in Appendix B.1.

4.2 Comparing CARSM and ARSM-MC

One major difference between CARSM and ARSM-MC is the usage of Q -Critic. It saves us from running

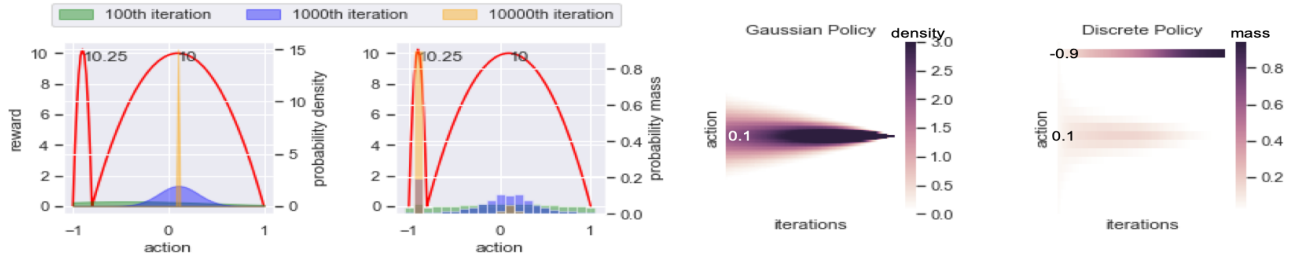


Figure 1: **left panel**: Change of policy over iterations in a single random trial between Gaussian policy (left) and discrete policy (right) on a bimodal-reward toy example. **right panel**: Average density on each action along with the training iteration between Gaussian and discrete policies for 100 random trials. Under this setting, the Gaussian policy fails to converge to the global optimum while discrete policy always finds the global optimum.

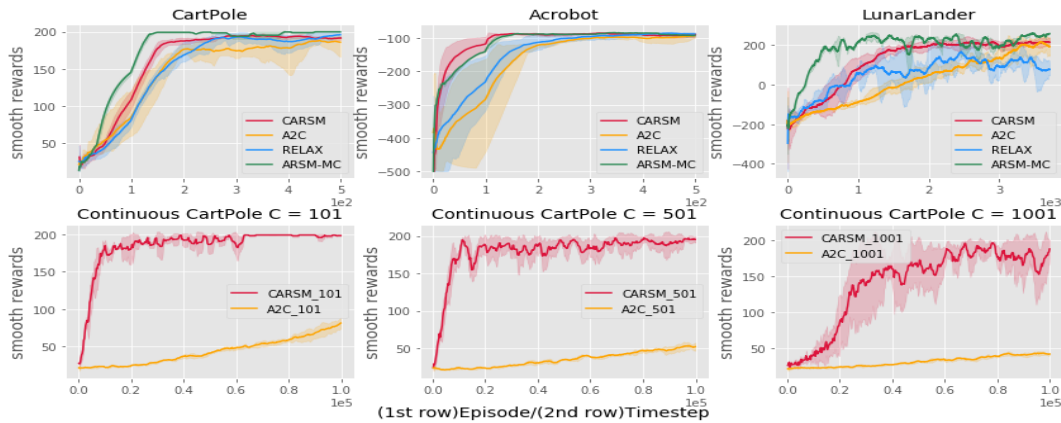


Figure 2: **top row**: Performance curves for discrete domains. Comparison between: A2C, RELAX, ARSM-MC, and CARSM. We show the cumulative rewards during training, moving averaged across 100 epochs; the curves show the mean \pm std performance across 5 random seeds. **bottom row**: Performance curves on CartPole with very large discrete action space. Comparison between: A2C and CARSM over a range of different discretization scale $C \in \{101, 501, 1001\}$. We show the cumulative rewards during training, moving averaged across 100 epochs; the curves show the mean \pm std performance across 5 random seeds.

MC rollouts to estimate the action-value functions of all unique pseudo actions, the number of which can be enormous under a multidimensional setting. This saving is at the expense of introducing bias to gradient estimation (not by the gradient estimator per se but by how Q is estimated). Similar to the argument between MC and TD, there is a trade-off between bias and variance. In this set of experiments, we show that the use of Critic in CARSM not only brings us accelerated training, but also helps return good performance.

To make the results of CARSM directly comparable with those of ARSM-MC shown in Yin et al. [2019], we evaluate the performances on an *Episode* basis on discrete classical-control tasks: CartPole, Acrobot, and LunarLander. We follow Yin et al. [2019] to limit the MC rollout sizes for ARSM-MC as 16, 64, and 1024, respectively. From Figure 2 top row, ARSM-MC has a better performance than CARSM on both CartPole and LunarLander, while CARSM outperforms the rest on Acrobot. The results are promising in the sense that CARSM only uses one rollout for estimation while ARSM-MC uses up to 16, 64, and 1024, respectively,

so CARSM largely improves the sample efficiency of ARSM-MC while maintaining comparable performance. The action space is uni-dimensional with 2, 3, and 4 discrete actions for CartPole, Acrobot, and LunarLander, respectively. We also compare CARSM and ARSM-MC given fixed number of timesteps. Under this setting, CARSM outperforms ARSM-MC by a large margin on both Acrobot and LunarLander. See Figure 7 in Appendix B.3 for more details.

4.3 Large Discrete Action Space

We want to show that CARSM has better sample efficiency on cases where the number of action C in one dimension is large. We test CARSM along with A2C on a continuous CartPole task, which is a modified version of discrete CartPole. In this continuous environment, we restrict the action space to $[-1, 1]$. Here the action indicates the force applied to the Cart at any time.

The intuition of why CARSM is expected to perform well under a large action space setting is because of the low-variance property. When C is large, the dis-

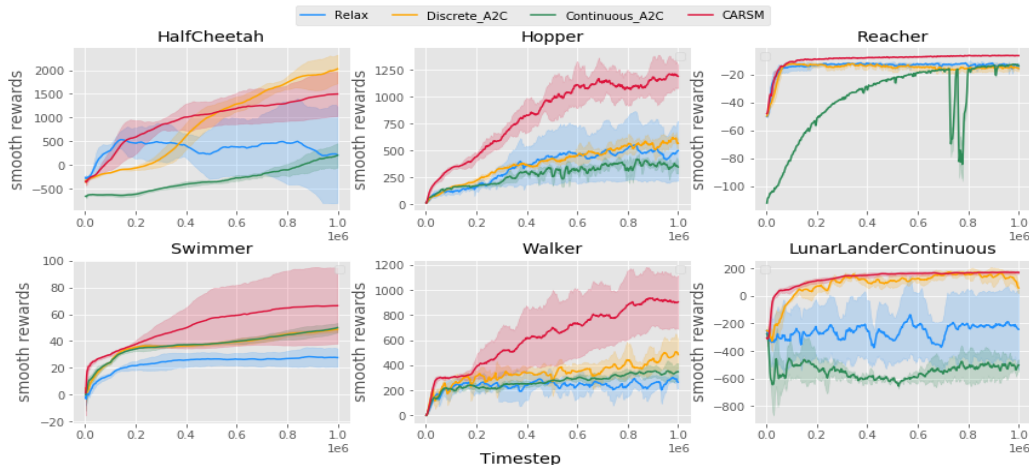


Figure 3: Performance curves on six benchmark tasks (all except the last are MuJoCo tasks). Comparison between: continuous A2C (Gaussian policy), discrete A2C, RELAX, and CARSM policy gradient. We show the cumulative rewards during training, moving averaged across 100 epochs; The curves show the mean \pm std performance across 5 random seeds.

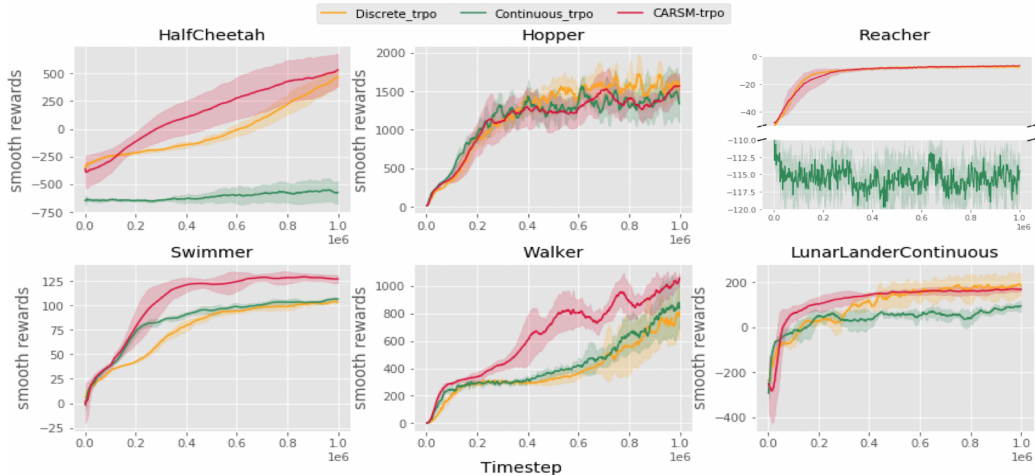


Figure 4: Performance curves on six benchmark tasks (all except the last are MuJoCo tasks). Comparison between: continuous TRPO (Gaussian policy), discrete TRPO, and CARSM policy gradient combined with TRPO. We show the cumulative rewards during training, moving averaged across 100 epochs; the curves show the mean \pm std performance across 5 random seeds.

tribution is more dispersed on each action compared with smaller case, which requires the algorithm captures the signal from best action accurately to improve exploitation. In this case, a high-variance gradient estimator will surpass the right signal, leading to a long exploration period or even divergence.

As shown in Figure 2 bottom row, CARSM outperforms A2C by a large margin in all three large C settings. Though the CARSM curve exhibits larger variations as C increases, it always learns much more rapidly at an early stage compared with A2C. Note the naive ARSM-MC algorithm will not work on this setting simply because it needs to run as many as tens of thousands MC rollouts to get a single gradient estimate.

4.4 OpenAI Gym Benchmark Tasks

In this set of experiments, we compare CARSM with A2C and RELAX, which all share the same underlying idea of improving the sample efficiency by reducing the variance of gradient estimation. For A2C, we compare with both Gaussian and discrete policies to check the intuition presented in Section 4.1. In all these tasks, following the results from Tang and Agrawal [2019], the action space is equally divided into $C = 11$ discrete actions at each dimension. Thus the discrete action space size becomes 11^K , where K is the action-space dimension that is 6 for HalfCheetah, 3 Hopper, 2 Reacher, 2 Swimmer, 6 Walker2D, and 2 LunarLander. More details on Appendix B.2.

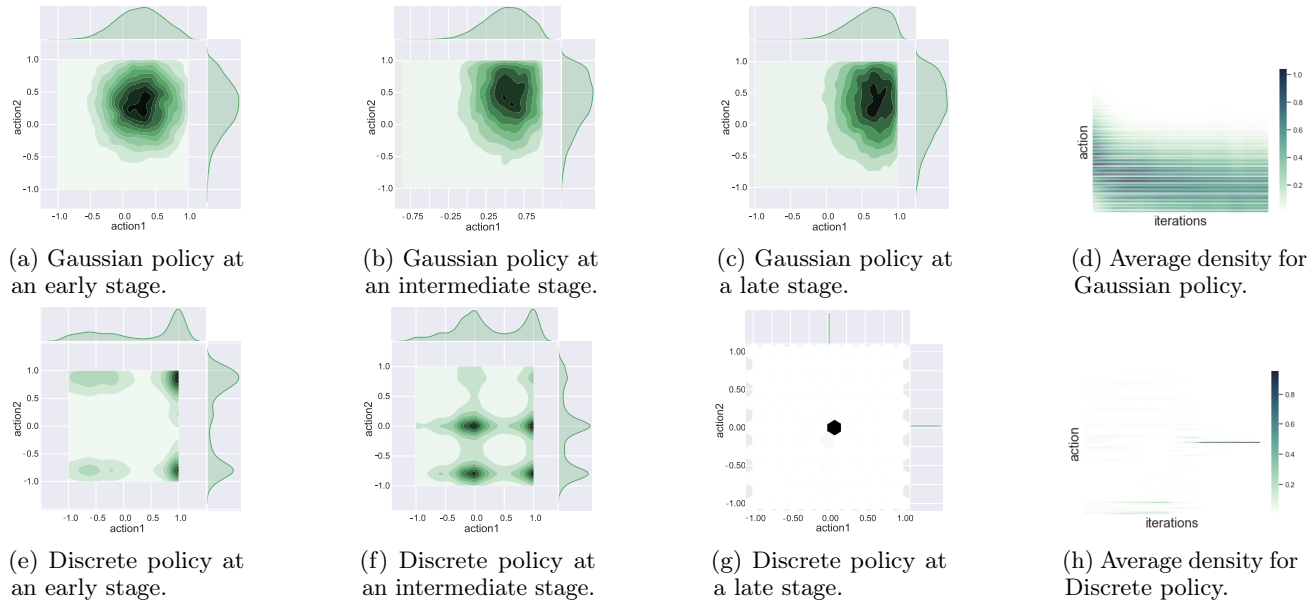


Figure 5: Policy distribution on the Reacher task between discrete policy and Gaussian policy for a given state (discrete action space has 11 actions on each dimension).

As shown in Figure 3, CARSM outperforms the other algorithms by a large margin except on HalfCheetah, demonstrating the high-sample efficiency of CARSM. Moreover, the distinct behaviors of Gaussian and discrete policies in the Reacher task, as shown in both Figures 3 and 4, are worth thinking, motivating us to go deeper on this task to search for possible explanations. We manually select a state that requires exploration on the early stage, and record the policy evolution along with training process at that specific state. We show those transition phases in Figure 5 for both Gaussian policy (top row) and discrete policy (bottom row).

For discrete policy, plots (e)-(g) in Figure 5 bottom row show interesting property: at the early stage, the policy does not put heavy mass at all on the final sub-optimal point $(0, 0)$, but explores around multiple density modes; then it gradually concentrates on several sup-optimal points on an intermediate phase, and converges to the final sub-optimal point. Plot (h) also conveys the same message that during the training process, discrete action can transit explorations around several density modes since the green lines can jump along the iterations. (The heatmaps of (d) and (h) in Figure 5 are computed in the same way as that in Figure 1, and details can be found in Appendix B.1.)

By contrast, Gaussian policy does not have the flexibility of exploring based on different density modes, therefore from plots (a)-(c) on the top row of Figure 5, the policy moves with a large radius but one center, and on (d), the green lines move consecutively which indicates a smooth but potentially not comprehensive exploration.

4.5 Combining CARSM with TRPO

Below we show that CARSM can be readily applied under TRPO to improve its performance. In the update step of TRPO shown in (4), the default estimator for $\nabla_{\theta} J(\theta)$ is A2C or its variant. We replace it with CARSM estimator and run it on the same set of tasks. As shown in Figure 4, Gaussian policy fails to find a good sub-optimal solution under TRPO for both HalfCheetah and Reacher and performs similarly to its discrete counterpart on the other tasks. Meanwhile, CARSM improves the performance of TRPO over discrete policy setting on three tasks and maintains similar performance on the others, which shows evidence that CARSM is an easy plug-in estimator for $\nabla_{\theta} J(\theta)$ and hence can potentially improve other algorithms, such as some off-policy ones [Wang et al., 2016, Degrís et al., 2012], that need this gradient estimation.

5 CONCLUSION

To solve RL tasks with multidimensional discrete action setting efficiently, we propose Critic-ARSM policy gradient, which is a combination of multidimensional sparse ARSM gradient estimator and an action-value critic, to improve sample efficiency for on-policy algorithm. We show the good performances of this algorithm from perspectives including stability on very large action space cases and comparisons with other standard benchmark algorithms, and show its potential to be combined with other standard algorithms. Moreover, we demonstrate the potential benefits of discretizing continuous control tasks to obtain a better exploration based on multimodal property.

Acknowledgements

This research was supported in part by Award IIS-1812699 from the U.S. National Science Foundation and the 2018-2019 McCombs Research Excellence Grant. The authors acknowledge the support of NVIDIA Corporation with the donation of the Titan Xp GPU used for this research, and the computational support of Texas Advanced Computing Center.

References

- M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018.
- S. Bhatnagar, D. Precup, D. Silver, R. S. Sutton, H. R. Maei, and C. Szepesvári. Convergent temporal-difference learning with arbitrary smooth function approximation. In *Advances in Neural Information Processing Systems*, pages 1204–1212, 2009.
- D. M. Blei, A. Kucukelbir, and J. D. McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- T. Degris, M. White, and R. S. Sutton. Off-policy actor-critic. *arXiv preprint arXiv:1205.4839*, 2012.
- G. Dulac-Arnold, R. Evans, H. van Hasselt, P. Sunehag, T. Lillicrap, J. Hunt, T. Mann, T. Weber, T. Degris, and B. Coppin. Deep reinforcement learning in large discrete action spaces. *arXiv preprint arXiv:1512.07679*, 2015.
- S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- W. Grathwohl, D. Choi, Y. Wu, G. Roeder, and D. Duvenaud. Backpropagation through the void: Optimizing control variates for black-box gradient estimation. *arXiv preprint arXiv:1711.00123*, 2017.
- S. Gu, T. Lillicrap, Z. Ghahramani, R. E. Turner, and S. Levine. Q-Prop: Sample-efficient policy gradient with an off-policy critic. *arXiv preprint arXiv:1611.02247*, 2016.
- T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1352–1361. JMLR. org, 2017.
- T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- Y. Hu, Q. Da, A. Zeng, Y. Yu, and Y. Xu. Reinforcement learning to rank in e-commerce search engine: Formalization, analysis, and application. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 368–377. ACM, 2018.
- E. Jang, S. Gu, and B. Poole. Categorical reparameterization with Gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- W. Jaśkowski, O. R. Lykkebø, N. E. Toklu, F. Triffterer, Z. Buk, J. Koutník, and F. Gomez. Reinforcement learning to run... fast. In *The NIPS'17 Competition: Building Intelligent Systems*, pages 155–167. Springer, 2018.
- D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- H. Liu, Y. Feng, Y. Mao, D. Zhou, J. Peng, and Q. Liu. Action-dependent control variates for policy optimization via Stein’s identity. *arXiv preprint arXiv:1710.11198*, 2017a.
- Y. Liu, P. Ramachandran, Q. Liu, and J. Peng. Stein variational policy gradient. *arXiv preprint arXiv:1704.02399*, 2017b.
- P. MacAlpine and P. Stone. UT Austin Villa: Robocup 2017 3d simulation league competition and technical challenges champions. In *Robot World Cup*, pages 473–485. Springer, 2017.
- H. R. Maei, C. Szepesvári, S. Bhatnagar, and R. S. Sutton. Toward off-policy learning control with function approximation. In *ICML*, pages 719–726, 2010.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing Atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937, 2016.

- OpenAI. OpenAI Five. <https://blog.openai.com/openai-five/>, 2018.
- T. Raiko, M. Berglund, G. Alain, and L. Dinh. Techniques for learning binary stochastic feedforward neural networks. *arXiv preprint arXiv:1406.2989*, 2014.
- J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015a.
- J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015b.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science*, 362(6419):1140–1144, 2018.
- R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- Y. Tang and S. Agrawal. Discretizing continuous action space for on-policy optimization. *arXiv preprint arXiv:1901.10500*, 2019.
- E. Todorov, T. Erez, and Y. Tassa. MuJoCo: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- G. Tucker, A. Mnih, C. J. Maddison, J. Lawson, and J. Sohl-Dickstein. REBAR: Low-variance, unbiased gradient estimates for discrete latent variable models. In *Advances in Neural Information Processing Systems*, pages 2627–2636, 2017.
- H. Van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double Q-learning. In *Thirtieth AAAI conference on artificial intelligence*, 2016.
- Z. Wang, V. Bapst, N. Heess, V. Mnih, R. Munos, K. Kavukcuoglu, and N. de Freitas. Sample efficient actor-critic with experience replay. *arXiv preprint arXiv:1611.01224*, 2016.
- C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. In *Reinforcement Learning*, pages 5–32. Springer, 1992.
- C. Wu, A. Rajeswaran, Y. Duan, V. Kumar, A. M. Bayen, S. Kakade, I. Mordatch, and P. Abbeel. Variance reduction for policy gradient with action-dependent factorized baselines. *arXiv preprint arXiv:1803.07246*, 2018.
- M. Yin and M. Zhou. ARM: Augment-REINFORCE-merge gradient for stochastic binary networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S11g0jAcYm>.
- M. Yin, Y. Yue, and M. Zhou. ARSM: Augment-REINFORCE-swap-merge estimator for gradient back-propagation through categorical variables. In *International Conference on Machine Learning*, 2019.

Discrete Action On-Policy Learning with Action-Value Critic: Supplementary Material

A Proof of Theorem 1

We first show the sparse ARSM for multidimensional action space case at one specific time point, then generalize it to stochastic setting. Since a_k are conditionally independent given ϕ_k , the gradient of ϕ_{kc} at one time point would be (we omit the subscript t for simplicity here)

$$\nabla_{\phi_{kc}} J(\phi) = \mathbb{E}_{\mathbf{a}_{\setminus k} \sim \prod_{k' \neq k} \text{Discrete}(a_{k'}; \sigma(\phi_{k'}))} [\nabla_{\phi_{kc}} \mathbb{E}_{a_k \sim \text{Cat}(\sigma(\phi_k))} [Q(\mathbf{a}, \mathbf{s})]],$$

and we apply the ARSM gradient estimator on the inner expectation part, which gives us

$$\begin{aligned} \nabla_{\phi_{kc}} J(\phi) &= \mathbb{E}_{\mathbf{a}_{\setminus k} \sim \prod_{k' \neq k} \text{Discrete}(a_{k'}; \sigma(\phi_{k'}))} \left\{ \mathbb{E}_{\boldsymbol{\varpi}_k \sim \text{Dir}(\mathbf{1}_C)} \left[(Q([\mathbf{a}_{\setminus k}, a_k^{c=j}], \mathbf{s}) - \frac{1}{C} \sum_{m=1}^C Q([\mathbf{a}_{\setminus k}, a_k^{m=j}], \mathbf{s})) (1 - C\varpi_{kj}) \right] \right\} \\ &= \mathbb{E}_{\boldsymbol{\varpi}_k \sim \text{Dir}(\mathbf{1}_C)} \left\{ \mathbb{E}_{\mathbf{a}_{\setminus k} \sim \prod_{k' \neq k} \text{Discrete}(a_{k'}; \sigma(\phi_{k'}))} \left[(Q([\mathbf{a}_{\setminus k}, a_k^{c=j}], \mathbf{s}) - \frac{1}{C} \sum_{m=1}^C Q([\mathbf{a}_{\setminus k}, a_k^{m=j}], \mathbf{s})) (1 - C\varpi_{kj}) \right] \right\} \end{aligned} \quad (6)$$

$$= \mathbb{E}_{\boldsymbol{\varpi}_k \sim \text{Dir}(\mathbf{1}_C)} \left\{ \mathbb{E}_{\prod_{k' \neq k} \text{Dir}(\boldsymbol{\varpi}_{k'}; \mathbf{1}_C)} \left[(Q(\mathbf{a}^{c=j}, \mathbf{s}) - \frac{1}{C} \sum_{m=1}^C Q(\mathbf{a}^{m=j}, \mathbf{s})) (1 - C\varpi_{kj}) \right] \right\}, \quad (7)$$

where (6) is derived by changing the order of two expectations and (7) can be derived by following the proof of Proposition 5 in Yin et al. [2019]. Therefore, if given $\boldsymbol{\varpi}_k \sim \text{Dir}(\mathbf{1}_C)$, it is true that $a_k^{c=j} = a_k$ for all (c, j) pairs, then the inner expectation term in (6) will be zero and consequently we have

$$g_{kc} = 0$$

as an unbiased single sample estimate of $\nabla_{\phi_{kc}} J(\phi)$; If given $\boldsymbol{\varpi}_k \sim \text{Dir}(\mathbf{1}_C)$, there exist (c, j) that $a_k^{c=j} \neq a_k$, we can use (7) to provide

$$g_{kc} = \sum_{j=1}^C \left[Q(\mathbf{s}, \mathbf{a}^{c=j}) - \frac{1}{C} \sum_{m=1}^C Q(\mathbf{s}, \mathbf{a}^{m=j}) \right] \left(\frac{1}{C} - \varpi_{kj} \right) \quad (8)$$

as an unbiased single sample estimate of $\nabla_{\phi_{kc}} J(\phi)$.

For a specific time point t , the objective function can be decomposed as

$$\begin{aligned} J(\phi_{0:\infty}) &= \mathbb{E}_{\mathcal{P}(\mathbf{s}_0)} [\prod_{t'=0}^{t-1} \mathcal{P}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \text{Cat}(\mathbf{a}_{t'}; \sigma(\phi_{t'}))] \left\{ \mathbb{E}_{\mathbf{a}_t \sim \text{Cat}(\sigma(\phi_t))} \left[\sum_{t'=0}^{t-1} \gamma^{t'} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) + \gamma^t Q(\mathbf{s}_t, \mathbf{a}_t) \right] \right\} \\ &= \mathbb{E}_{\mathcal{P}(\mathbf{s}_0)} [\prod_{t'=0}^{t-1} \mathcal{P}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \text{Cat}(\mathbf{a}_{t'}; \sigma(\phi_{t'}))] \left\{ \mathbb{E}_{\mathbf{a}_t \sim \text{Cat}(\sigma(\phi_t))} \left[\sum_{t'=0}^{t-1} \gamma^{t'} r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) \right] \right\} \\ &\quad + \mathbb{E}_{\mathcal{P}(\mathbf{s}_0)} [\prod_{t'=0}^{t-1} \mathcal{P}(\mathbf{s}_{t'+1} | \mathbf{s}_{t'}, \mathbf{a}_{t'}) \text{Cat}(\mathbf{a}_{t'}; \sigma(\phi_{t'}))] \left\{ \mathbb{E}_{\mathbf{a}_t \sim \text{Cat}(\sigma(\phi_t))} [\gamma^t Q(\mathbf{s}_t, \mathbf{a}_t)] \right\}, \end{aligned}$$

where the first part has nothing to do with ϕ_t , we therefore have

$$\nabla_{\phi_{tkc}} J(\phi_{0:\infty}) = \mathbb{E}_{\mathcal{P}(\mathbf{s}_t | \mathbf{s}_0, \pi_\theta) \mathcal{P}(\mathbf{s}_0)} \left\{ \gamma^t \nabla_{\phi_{tkc}} \mathbb{E}_{\mathbf{a}_t \sim \text{Cat}(\sigma(\phi_t))} [Q(\mathbf{s}_t, \mathbf{a}_t)] \right\}$$

With the result from (8), the statements in Theorem1 follow.

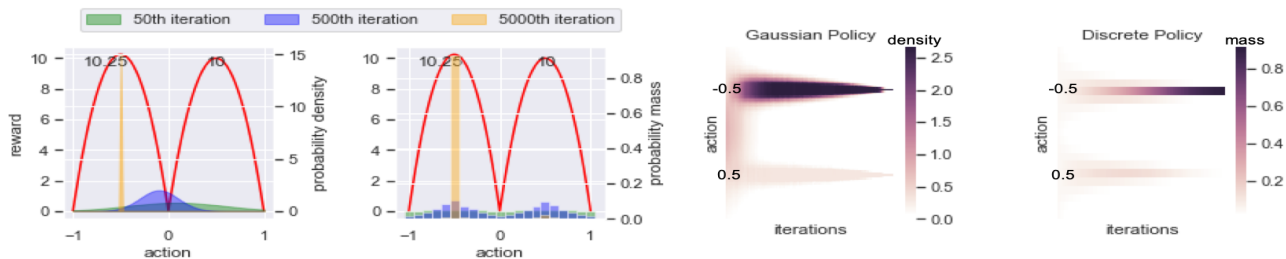


Figure 6: **left panel:** Change of policy over iterations between Gaussian policy (left) and discrete policy (right) on toy example setting. **right panel:** Average density on each action along with the training iterations between Gaussian policy and discrete policy for 100 experiments. (The Gaussian policy converges to the inferior optimal solution 12 times out of 100 times, and discrete policy converges to the global optimum all the time).

B Experiment setup

B.1 Toy example setup

Assume the true reward is a bi-modal distribution (as shown in Figure 6 left panel red curves) with a difference between its two peaks:

$$r(a) = \begin{cases} -c_1(a-1)(a-m) + \epsilon_1 & \text{for } a \in [m, 1] \\ -c_2(a+1)(a-m) + \epsilon_2 & \text{for } a \in [-1, m], \end{cases}$$

where the values of c_1 , c_2 , and m determine the heights and widths of these two peaks, and $\epsilon_1 \sim \mathcal{N}(0, 2)$ and $\epsilon_2 \sim \mathcal{N}(0, 1)$ are noise terms. It is clear that $a_{\text{left}}^* = (m-1)/2$ and $a_{\text{right}}^* = (1+m)/2$ are two local-optimal solutions and corresponding to $r_{\text{left}} := \mathbb{E}[r(a_{\text{left}}^*)] = c_2(1+m)^2/4$ and $r_{\text{right}} := \mathbb{E}[r(a_{\text{right}}^*)] = c_1(1-m)^2/4$. Here we always choose c_1 and c_2 such that r_{left} is slightly bigger than r_{right} which makes a_{left}^* a better local-optimal solution. It is clear that the more closer a_{left}^* to -1 , the more explorations a policy will need to converge to a_{left}^* . Moreover, the noise terms can give wrong signals and may lead to bad update directions, and exploration will play an essential role in preventing the algorithm from acting too greedily. The results shown on Section 4.1 has $m = -0.8$, $c_1 = 40/(1.8^2)$ and $c_2 = 41/(0.2^2)$, which makes $r_{\text{left}} = 10.25$ and $r_{\text{right}} = 10$. We also show a simple example at Figure 6 with $m = 0$, $c_1 = 40/(0.5^2)$ and $c_2 = 41/(0.5^2)$, which maintains the same peak values.

The experiment setting is as follows: for each episode, we collect 100 samples and update the corresponding parameters ($[\mu, \sigma]$ for Gaussian policy and $\phi \in \mathbb{R}^{21}$ for discrete policy where the action space is discretized to 21 actions), and iterate until N samples are collected. We add a quadratic decaying coefficient for the entropy term for both policies to encourage explorations on an early stage. The Gaussian policy is updated using reparametrization trick [Kingma and Welling, 2013], which can be applied to this example since we know the derivative of the reward function (note this is often not the case for RL tasks). The discrete policy is updated using ARSM gradient estimator described in Section 2.

On the heatmap, the horizontal axis is the iterations, and vertical axis denotes the actions. For each entry corresponding to a at iteration i , its value is calculated by $v(i, a) = \frac{1}{U} \sum_{u=1}^U p_u(a|i)$, where $p_u(a|i)$ is the probability of taking action a at iteration i for that policy in u th trial.

We run the same setting with different seeds for Gaussian policy and discrete policy for 100 times, where the initial parameters for Gaussian Policy is $\mu_0 = m, \sigma = 1$ and for discrete policy is $\phi_i = 0$ for any i to eliminate the effects of initialization.

In those 100 trials, when $m = -0.8, N = 1e^6$, Gaussian policy fails to find the true global optimal solution (0/100) while discrete policy can always find that optimal one (100/100). When $m = 0, N = 5e^5$, the setting is easier and Gaussian policy performs better in this case with only 12/100 percentage converging to the inferior sub-optimal point 0.5, and the rest 88/100 chances getting to global optimal solution. On the other hand, discrete policy always converges to the global optimum (100/100). The similar plots are shown on Figure 6. The p -value for this proportion test is 0.001056, which shows strong evidence that discrete policy outperforms Gaussian policy on this example.

B.2 Baselines and CARSM setup

Our experiments aim to answer the following questions: **(a)** How does the proposed CARSM algorithm perform when compared with ARSM-MC (when ARSM-MC is not too expensive to run). **(b)** Is CARSM able to efficiently solve tasks with large discrete action spaces (i.e., C is large). **(c)** Does CARSM have better sample efficiency than the algorithms, such as A2C and RELAX, that have the same idea of using baselines for variance reduction. **(d)** Can CARSM combined with other standard algorithms such as TRPO to achieve a better performance.

Baselines and Benchmark Tasks. We evaluate our algorithm on benchmark tasks on OpenAI Gym classic-control and MuJoCo tasks [Todorov et al., 2012]. We compare the proposed CARSM with ARSM-MC [Yin et al., 2019], A2C [Mnih et al., 2016], and RELAX [Grathwohl et al., 2017]; all of them rely on introducing baseline functions to reduce gradient variance, making it fair to compare them against each other. We then integrate CARSM into TRPO by replacing the A2C gradient estimator for $\nabla_{\theta} J(\theta)$, and evaluate the performances on MuJoCo tasks to show that a simple plug-in of the CARSM estimator can bring the improvement.

Hyper-parameters: Here we detail the hyper-parameter settings for all algorithms. Denote β_{policy} and β_{critic} as the learning rates for policy parameters and Q critic parameters, respectively, n_{critic} as the number of training time for Q critic, and α as the coefficient for entropy term. For CARSM, we select the best learning rates $\beta_{\text{policy}}, \beta_{\text{critic}} \in \{1, 3\} \times 10^{-2}$, and $n_{\text{critic}} \in \{50, 150\}$; For A2C and RELAX, we select the best learning rates $\beta_{\text{policy}} \in \{3, 30\} \times 10^{-5}$. In practice, the loss function consists of a policy loss L_{policy} and value function loss L_{value} . The policy/value function are optimized jointly by optimizing the aggregate objective at the same time $L = L_{\text{policy}} + cL_{\text{value}}$, where $c = 0.5$. Such joint optimization is popular in practice and might be helpful in cases where policy/value function share parameters. For A2C, we apply a batched optimization procedure: at iteration t , we collect data using a previous policy iterate π_{t-1} . The data is used for the construction of a differentiable loss function L . We then take v_{iter} gradient updates over the loss function objective to update the parameters, arriving at π_t . In practice, we set $v_{\text{iter}} = 10$. For TRPO and TRPO combined with CARSM, we use max KL-divergence of 0.01 all the time without tuning. All algorithms use a initial α of 0.01 and decrease α exponentially, and target network parameter τ is 0.01. To guarantee fair comparison, we only apply the tricks that are related to each algorithm and didn't use any general ones such as normalizing observation. More specifically, we replace Advantage function with normalized Generalized Advantage Estimation (GAE) [Schulman et al., 2015b] on A2C, apply normalized Advantage on RELAX.

Structure of Q critic networks: There are two common ways to construct a Q network. The first one is to model the network as $Q : \mathbb{R}^{n_S} \rightarrow \mathbb{R}^{|\mathcal{A}|}$, where n_S is the state dimension and $|\mathcal{A}| = C^K$ is the number of unique actions. The other structure is $Q : \mathbb{R}^{n_S+K} \rightarrow \mathbb{R}$, which means we need to concatenate the state vector \mathbf{s} with action vector \mathbf{a} and feed that into the network. The advantage of first structure is that it doesn't involve the issue that action vector and state vector are different in terms of scale, which may slow down the learning process or make it unstable. However, the first option is not feasible under most multidimensional discrete action situations because the number of actions grow exponentially along with the number of dimension K . Therefore, we apply the second kind of structure for Q network, and update Q network multiple times before using it to obtain the CARSM estimator to stabilize the learning process.

Structure of policy network: The policy network will be a function of $\mathcal{T}_{\theta} : \mathbb{R}^{n_S} \rightarrow \mathbb{R}^{K \times C}$, which feed in state vector \mathbf{s} and generate $K \times C$ logits ϕ_{kc} . Then the action is obtained for each dimension k by $\pi(a_k | \mathbf{s}, \theta) = \sigma(\phi_k)$, where $\phi_k = (\phi_{k1}, \dots, \phi_{kC})'$. For both the policy and Q critic networks, we use a two-hidden-layer multilayer perceptron with 64 nodes per layer and tanh activation.

Environment setup

- *HalfCheetah* ($\mathcal{S} \subset \mathbb{R}^{17}, \mathcal{A} \subset \mathbb{R}^6$)
- *Hopper* ($\mathcal{S} \subset \mathbb{R}^{11}, \mathcal{A} \subset \mathbb{R}^3$)
- *Reacher* ($\mathcal{S} \subset \mathbb{R}^{11}, \mathcal{A} \subset \mathbb{R}^2$)
- *Swimmer* ($\mathcal{S} \subset \mathbb{R}^8, \mathcal{A} \subset \mathbb{R}^2$)
- *Walker2D* ($\mathcal{S} \subset \mathbb{R}^{17}, \mathcal{A} \subset \mathbb{R}^6$)
- *LunarLander Continuous* ($\mathcal{S} \subset \mathbb{R}^8, \mathcal{A} \subset \mathbb{R}^2$)

B.3 Comparison between CARSM and ARSM-MC for fixed timestep

We compare ARSM-MC and CARSM for fixed timestep setting, with their performances shown in Figure 7

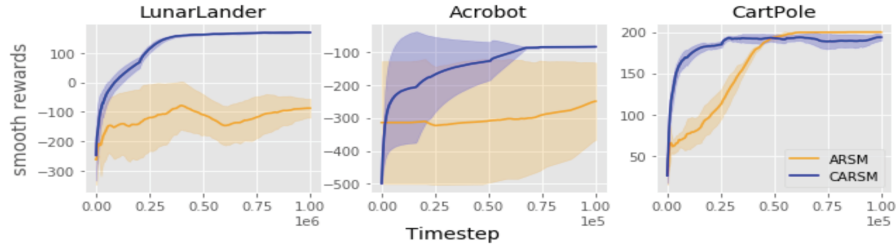


Figure 7: Performance curves for comparison between ARSM-MC and CARSM given fix timesteps

C Pseudo Code

We provide detailed pseudo code to help understand the implementation of CARSM policy gradient. There are four major steps for each update iteration: (1) Collecting samples using augmented Dirichlet variables ϖ_t ; (2) Update the Q critic network using both on-policy samples and off-policy samples; (3) Calculating the CARSM gradient estimator; (4) soft updating the target networks for both the policy and critic. The (1) and (3) steps are different from other existing algorithms and we show their pseudo codes in Algorithms 1 and 2, respectively.

Algorithm 1: Collecting samples from environment

Input: Policy network $\pi(\mathbf{a} | \mathbf{s}, \boldsymbol{\theta})$, initial state \mathbf{s}_0 , sampled step T , replay buffer \mathcal{R}

Output: Intermediate variable matrix $\varpi_{1:T}$, logit variables $\phi_{1:T}$, rewards vector $r_{1:T}$, state vectors $\mathbf{s}_{1:T}$, action vectors $\mathbf{a}_{1:T}$, replay buffer \mathcal{R}

for $t = 1 \dots T$ **do**

 Generate Dirichlet random variable $\varpi_{tk} \sim \text{Dir}(\mathbf{1}_C)$ for each dimension k ;

 Calculate logits $\phi_t = \mathcal{T}_{\boldsymbol{\theta}}(\mathbf{s}_t)$ which is a $K \times C$ length vector

 Select action $a_{tk} = \text{argmin}_{i \in \{1, \dots, C\}} (\ln \varpi_{tki} - \phi_{tki})$ for each dimension k ;

 Obtain next state values \mathbf{s}_{t+1} and reward r_t based actions $\mathbf{a}_t = (a_{t1}, \dots, a_{tK})'$ and current state \mathbf{s}_t .

 Store the transition $\{\mathbf{s}_t, \mathbf{a}_t, r_t, \mathbf{s}_{t+1}\}$ to replay buffer \mathcal{R}

 Assign $\mathbf{s}_t \leftarrow \mathbf{s}_{t+1}$.

end for

Algorithm 2: CARSM policy gradient for a K -dimensional C -way categorical action space.

Input: Critic network Q_ω , policy network π_θ , on-policy samples including states $\mathbf{s}_{1:T}$, actions $\mathbf{a}_{1:T}$, intermediate Dirichlet random variables $\varpi_{1:T}$, logits vectors $\phi_{1:T}$, discounted cumulative rewards $y_{1:T}$.

Output: an updated policy network

Initialize $g \in \mathbb{R}^{T \times K \times C}$;

for $t = 1 \dots T$ (in parallel) **do**

for $k = 1 \dots K$ (in parallel) **do**

 Let $A_{tk} = \{(c, j)\}_{c=1:C, j < c}$, and initialize $P^{tk} \in \mathbb{R}^{C \times C}$ with all element equals to a_{tk} (true action).

for $(c, j) \in A_{tk}$ (in parallel) **do**

 Let $a_{tk}^{c=j} = \arg \min_{i \in \{1, \dots, C_k\}} (\ln \varpi_{tki}^{c=j} - \phi_{tki})$

if $a_{tk}^{c=j}$ not equals to a_{tk} **then**

 | Assign $a_{tk}^{c=j}$ to $P^{tk}(c, j)$

end if

end for

end for

Let $S_t = \text{unique}(P^{t1} \otimes P^{t2} \dots \otimes P^{tK}) \setminus \{a_{t1} \otimes a_{t2} \dots \otimes a_{tK}\}$, which means S_t is the set of all unique values across K dimensions except for true action $\mathbf{a}_t = \{a_{t1} \otimes a_{t2} \dots \otimes a_{tK}\}$; denote pseudo action of swapping between coordinate c and j as $S_t(c, j) = (P^{t1}(c, j) \otimes P^{t1}(c, j) \dots \otimes P^{tK}(c, j))$, and define \mathcal{I}_t as unique pairs contained in S_t .

Initialize matrix $F^t \in \mathbb{R}^{C \times C}$ with all elements equal to y_t ;

for $(\tilde{c}, \tilde{j}) \in \mathcal{I}_t$ (in parallel) **do**

 | $F^t(\tilde{c}, \tilde{j}) = Q_\omega(\mathbf{s}_t, S_t(\tilde{c}, \tilde{j}))$

end for

Plug in number for matrix $g_{tkc} = \sum_{j=1}^C (F_c^t - \bar{F}_c^t)(\frac{1}{C} - \varpi_{tkj})$, where F_c^{tk} denotes the c th row of matrix F^t and \bar{F}_c^t is the mean of that row;

for $k = 1 \dots K$ **do**

if every element in P^{tk} is a_{tk} **then**

 | $g_{tkc} = 0$

end if

end for

end for

Update the parameter for θ for policy network by maximize the function

$$J = \frac{1}{TKC} \sum_{t=1}^T \sum_{k=1}^K \sum_{c=1}^C g_{tkc} \phi_{tkc}$$

where ϕ_{tkc} are logits and g_{tkc} are placeholders that stop any gradients, and use auto-differentiation on ϕ_{tkc} to obtain gradient with respect to θ .
