# ON HIGH-DIMENSIONAL ACTION SELECTION FOR DEEP REINFORCEMENT LEARNING

### Anonymous authors Paper under double-blind review

### ABSTRACT

With recent advances in deep reinforcement learning (RL), *high-dimensional action selection* has become an important yet challenging problem in many real applications, especially in unknown and complex environments. Existing works often require a sophisticated prior design to eliminate redundancy in the action space, relying heavily on domain expert experience or involving high computational complexity, which limits their generalizability across different RL tasks. In this paper, we address these challenges by proposing a general data-driven action selection approach with model-free and computational-friendly properties. Our method not only *selects minimal sufficient actions* but also *controls the false discovery rate* via knockoff sampling. More importantly, we seamlessly integrate the action selection into deep RL methods during online training. Empirical experiments validate the established theoretical guarantees, demonstrating that our method surpasses various alternative techniques in terms of both performances in variable selection and overall achieved rewards.

### 1 INTRODUCTION

027 028

004 005

006

007 008 009

010 011

012

013

014

015

016

017

018

019

021

023

025 026

029 Recent advances in deep reinforcement learning (RL) have attracted significant attention, with applications spanning numerous fields such as robotics, games, healthcare, and finance [23; 21; 24; 54]. Despite their ability to handle sequential decision-making, the practical utility of RL methods in 031 real-world scenarios is often limited, especially in dealing with the high-dimensional action spaces [48; 21; 42; 53]. High-dimensional action spaces are prevalent in "black box" systems, characterized 033 by overloaded actionable variables that are often *abundant and redundant*. Examples include precision 034 medicine, where numerous combinations of treatments and dosages are possible [see e.g., 19; 31]; neuroscience, which involves various stimulation points and intensities [see e.g., 12]; and robotics, particularly in muscle-driven robot control, where coordination of numerous muscles is required [see 037 e.g., 44]. Nevertheless, these high-dimensional action spaces often contain many actions that are 038 either ineffective or have negligible impact on states and rewards. Training RL models on the entire action space can result in substantial inefficiencies in both computation and data collection.

040 To handle high dimensionality, a promising approach is to employ automatic dimension reduction 041 techniques to reduce the size of the action space, by selecting only the essential minimum action set 042 necessary for effectively learning the environment and optimizing the policy based on the subspace. 043 Having such a minimal yet sufficient action space can significantly enhance learning efficiency, as 044 agents can thoroughly explore a more concise set of actions [55; 22; 17]. Moreover, a smaller action space can reduce computational complexity, a notable benefit in deep RL, where neural networks are used for function approximation [47; 41]. In practical scenarios, eliminating superfluous actions saves 046 the cost of extensive measurement equipment and thus allows a more comprehensive exploration 047 of available actions. Yet, existing works often require a sophisticated prior design to eliminate 048 redundancy in the action space [e.g., 49; 18; 11; 32], relying heavily on domain expert experience or involving high computational complexity, limiting their generalizability across different RL tasks. 050

In this paper, we propose a general data-driven action selection approach to identify the minimum sufficient actions in the high-dimensional action space. To handle the complex environments often seen in deep RL, we develop a novel variable selection approach called knockoff sampling (KS) for online RL, with theoretical guarantees of *false discovery rate control*, inspired by the model-free

056

060

061

062

063

077

078

079

081

082

083

084

085

087

088

090



Figure 1: Average rewards under three proximal policy optimization (PPO) methods in a synthetic environment with 54 actions (among which only 4 actions influence states and rewards). The green line refers to the PPO trained based on the true influential actions, the red refers to the PPO with the estimated minimal sufficient actions by the proposed variable selection (KS), the black line represents the PPO with the entire redundant action space, and the dashed line is the optimal reward. The red line outperforms the black line indicating *the effectiveness of the variable selection* step.

knockoff method [7]. The effectiveness of this action selection method is demonstrated in Fig. 1.
Here, a proximal policy optimization (PPO) method [43] using variable selection outperforms the
one trained on the entire action space and shows comparable performance to the PPO trained with the
pre-known true minimal sufficient action. To remain computational-friendly, we design an adaptive
strategy with a simple mask operation that seamlessly integrates this action selection method into
deep RL methods during online training. Our main contributions are fourfold:

• Conceptually, this work pioneers exploring high-dimensional action selection in online RL. We formally define *the sufficient action set* encompassing all influential actions, and *the minimal sufficient action set*, which contains the smallest number of actions necessary for effective decision-making.

• Methodologically, our method *bypasses the common challenge of creating accurate knockoff features in model-free knockoffs*. We use the established distribution of actions from the current policy network in online RL to precisely resample action values, producing exact knockoff features.

• Algorithm-wise, to *flexibly integrate arbitrary variable selection into deep RL* and eliminate the need to initialize a new RL model after the selection, we design a binary hard mask approach based on the indices of selected actions. This efficiently neutralizes the influence of non-chosen actions.

• Theoretically, to *address the issues of highly dependent data* in online RL, we couple our KS method with sample splitting and majority vote; under commonly imposed conditions, we theoretically show our method *consistently* identifies the minimal sufficient action set with false discovery rate control.

091 1.1 RELATED WORKS

092 Deep reinforcement learning has made significant breakthroughs in complex sequential decisionmaking across various tasks [36; 46; 43; 15]. Yet, several considerable obstacles exist when dealing 094 with high-dimensional spaces using deep RL. In terms of high dimensional state space, the state abstraction [35; 37] has been studied to learn a mapping from the original state space to a much 096 smaller abstract space to preserve the original Markov decision process. Yet, these methods such as bisimulation can be computationally expensive and challenging when the state space is very large or 098 has complex dynamics [40]. Tied to our topic, it is hard to utilize such abstraction-based methods to implement transformed actions. This redirects us to variable selection on the redundant state space [see e.g., 25; 14]. Recently, Hao et al. [16] combined LASSO with fitted Q-iteration to reduce 100 states; following this context, Ma et al. [34] employs the knockoff method for state selection but with 101 discrete action spaces. However, all these works focus on the high dimensional state space in offline 102 data, while our method aims to extract sufficient and necessary actions during online learning. 103

For RL with the high-dimensional action space, especially for continuous actions, some studies
 [11; 49] transformed the continuous control problem into the combinatorial action problem, by
 discretizing large action spaces into smaller subspaces. However, this transformation can lead to a
 significant loss of precision and hence produce suboptimal solutions [27; 50]. Other works [see e.g.,
 18; 32] focused on muscle control tasks and used architectures reducing the action dimensionality

before deploying RL methods. One recent study by Schumacher et al. [44] combined differential
extrinsic plasticity with RL to control high-dimensional large systems. Yet, all these works require
specialized data collection, known joint ranges of actions, forced dynamics, or desired behaviors
of policies, before implementing RL. In contrast, our method is entirely *data-driven without prior knowledge of environments*, and thus can be generalized to tasks beyond muscle control.

113 Variable selection, also known as feature selection, is a critical process to choose the most relevant 114 variables representing the target outcome of interest, enhancing both model performance and interpre-115 tation. Over the past few decades, many well-known methods have been established, ranging from 116 classical LASSO, Fisher score, and kernel dimension reduction [51; 13; 8], towards deep learning 117 [3; 28; 26]. Yet, these works either suffer from model-based constraints or lack theoretical guarantees. The model-X knockoff method proposed by Candes et al. [7] aims to achieve both goals via a 118 general variable selection framework for black-box algorithms with guarantees of false discovery rate 119 control. Due to its model-agnostic nature, the knockoff method has been extended to complement a 120 wide range of variable selection approaches [33; 45; 29]. The main price or central challenge within 121 the knockoff method lies in the generation of faithful knockoff features. Existing techniques either 122 using model-specific methods [see e.g., 45; 30] that assume the underlying covariate distribution, 123 or model-free approaches [see e.g., 20; 39] that utilize deep generative models to obtain knockoffs 124 without further assumptions on feature distribution. Owing to the blessing of online RL, our method 125 bypasses this challenge through the known joint distribution of actions represented by the ongoing 126 policy network, and thus can easily resample the action values to create exact knockoff features.

127 128 129

130

131

### 2 PROBLEM SETUP

### 2.1 NOTATIONS

132 Consider a Markov Decision Process (MDP) characterized by the tuple  $(S, A, p, r, \gamma)$ , in which both 133 the state space S and the action space A are continuous. The state transition probability, denoted as 134  $p: \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0, \infty)$ , is an unknown probability density function that determines the likelihood 135 of transitioning to a next state  $\mathbf{s}_{t+1} \in \mathcal{S}$ , given the current state  $\mathbf{s}_t \in \mathcal{S}$  and the action  $\mathbf{a}_t \in \mathcal{A}$ . 136 The environment provides a reward, bounded within  $[r_{\min}, r_{\max}]$ , for each transition, expressed as 137  $r: S \times A \rightarrow [r_{\min}, r_{\max}]$ . The discount factor, represented by  $\gamma \in (0, 1)$ , influences the weighting of future rewards. We denote a generic tuple consisting of the current state, action, reward, and 138 subsequent state as  $(\mathbf{S}_t, \mathbf{A}_t, \mathbf{R}_t, \mathbf{S}_{t+1})$ . The Markovian property of MDP is that given the current 139 state  $S_t$  and action  $A_t$ , the current  $R_t$  and the next state  $S_{t+1}$ , are conditionally independent of the 140 past trajectory history. Consider  $\mathbf{A}_t \in \mathbb{R}^p$  where p is very large indicating a high dimensional action 141 space. We utilize  $\rho_{\pi}(\mathbf{s}_t)$  and  $\rho_{\pi}(\mathbf{s}_t, \mathbf{a}_t)$  to denote the state and state-action marginal distributions, 142 respectively, of the trajectory distribution generated by a policy  $\pi(\mathbf{a}_t \mid \mathbf{s}_t)$ . The notation  $J(\pi)$  is used to represent the expected discounted reward under this policy:  $J(\pi) = \sum_t \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_{\pi}} \left[ \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$ . 143 144 The goal of RL is to maximize the expected sum of discounted rewards above. This can be extended 145 to a more general maximum entropy objective with the expected entropy of the policy over  $\rho_{\pi}(\mathbf{s}_t)$ .

146 147

148

153 154

155

### 2.2 MINIMAL SUFFICIENT ACTION SET IN ONLINE RL

To address the high dimensional action space, we propose to utilize the variable selection instead of representation for practical usefulness. To achieve this goal, we first formally define the minimal sufficient action set. Denote the subvector of  $\mathbf{A}_t$  indexed by components in G as  $\mathbf{A}_{t,G}$  with an index set  $G \subseteq \{1, 2, ..., p\}$ . Let  $G^c = \{1, ..., p\} \setminus G$  be the complement of G.

**Definition 2.1.** (Sufficient Action Set) We say G is the sufficient action (index) set in an MDP if

 $R_t \perp \mathbf{A}_{t,G^c} \mid \mathbf{S}_t, \mathbf{A}_{t,G}, \qquad \mathbf{S}_{t+1} \perp \mathbf{A}_{t,G^c} \mid \mathbf{S}_t, \mathbf{A}_{t,G}, \qquad \text{for all } t \ge 0.$ 

The sufficient action set can be seen as a sufficient conditional set to achieve past and future independence. The sufficient action set may not be unique.

**Definition 2.2.** (Minimal Sufficient Action Set) We say G is the minimal sufficient action set in an MDP if it has the smallest cardinality among all sufficient action sets.

160

Unlike the sufficient action set, there is only one unique minimal sufficient action set to achieve

161 Unlike the sufficient action set, there is only one unique minimal sufficient action set to achieve conditional independence if there are no identical action variables in the environment. We also call

162  $G^{c}$  the redundant set when G is the minimal sufficient action (index) set. Here, to achieve such a 163 minimal sufficient action set, one should also require the states  $S_t$  be the sufficient states, so there is 164 no useless state [34] to introduce related redundant actions that possibly lead to ineffective exploration 165 or data inefficiency. Without loss of generality, we assume sufficient states throughout this paper and 166 focus on eliminating redundant actions in a high-dimensional action space. Our goal is to identify the minimal sufficient action set for online deep reinforcement learning to improve exploration. 167

168 169

170

181

### 2.3 PRELIMINARY: KNOCKOFF VARIABLE SELECTION

Without making additional assumptions on the dependence among variables, in this work, we utilize 171 the model-X knockoffs [7] for flexible variable selection, which ensures finite-sample control of 172 the false discovery rate (FDR). We first briefly review the model-X knockoffs [7] in the supervised 173 regression setting with independent samples, which will be leveraged later as the base variable 174 selector of our proposed method for dependent data in online RL setting. Specifically, given n175 independent observations, consider Y as a n-dimensional response vector and  $\mathbf{X} = (\mathbf{x}_1, \cdots, \mathbf{x}_p)$ 176 as an  $n \times p$  matrix of covariates. The knockoff inference aims to identify significant covariates that 177 influence the outcome while controlling FDR. Towards this goal, the model-X knockoff generates an 178  $n \times p$  matrix  $\mathbf{X} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_p)$  as knockoff features that have the similar properties as the collected 179 covariates. This matrix is constructed by the joint distribution of X and satisfy: 180

$$\widetilde{\mathbf{X}} \perp \mathbf{Y} \mid \mathbf{X} \quad \text{and} \quad (\mathbf{X}, \widetilde{\mathbf{X}})_{\operatorname{swap}(\Omega)} \stackrel{d}{=} (\mathbf{X}, \widetilde{\mathbf{X}}),$$
(1)

182 for each subset  $\Omega$  within the set  $\{1, \dots, p\}$ , where swap $(\Omega)$  indicates the operation of swapping such 183 that for each  $j \in \Omega$ , the j-th and (j + p)-th columns are interchanged. The notation  $\stackrel{d}{=}$  signifies 184 equality in distribution. After obtaining knockoff features, let  $\widetilde{\mathcal{D}} = \{X, \widetilde{X}, Y\}$  denote an augmented 185 dataset and we can calculate the feature importance scores  $Z_j$  and  $\widetilde{Z}_j$  for each variable  $x_j$  and its 186 corresponding knockoff  $\widetilde{x}_j$ . Define the function  $f: \mathbb{R}^2 \to \mathbb{R}$  as an anti-symmetric function, meaning 187 that f(u,v) = -f(v,u) for all  $u, v \in \mathbb{R}^2$ , e.g., f(u,v) = u - v. Set  $W_j = f(Z_j, Z_j)$  in such a 188 189 way that higher values of  $W_i$  indicate stronger evidence of the significance of  $x_i$  being influential 190 covariate. The j-th variable is selected if its corresponding  $W_i$  is at least a certain threshold  $\tau_{\alpha}$ 191 when the target FDR level is  $\alpha$ . For example, the set of chosen variables can be represented as  $\mathcal{I} = \{j : W_j \ge \tau_\alpha\}, \text{ where }$ 192

$$\tau_{\alpha} = \min\left\{\tau > 0: \frac{\#\{j \in [p]: W_j \le -\tau\}}{\#\{j \in [p]: W_j \ge \tau\}} \le \alpha\right\}.$$
(2)

195 196 197

199

201

203

193 194

#### 3 **ONLINE DEEP RL WITH VARIABLE SELECTION**

To identify the minimal sufficient action set in online deep RL, we integrate the action selection into RL to find truly influential actions during the training process. Its advantages are manifold. 200 Firstly, its model-agnostic nature ensures compatibility across various RL architectures and algorithms. Moreover, its data-driven characteristic allows for straightforward application across diverse scenarios, 202 thereby increasing practical utility. Crucially, the action selection boosts the explainability and reliability of RL systems by clearly delineating actions that contribute to model performance. In the 204 following, we first introduce an action-selected exploration strategy for online deep RL in Section 3.1, followed by the model-free knockoff-sampling method for action selection in Section 3.2.

### 3.1 ACTION-SELECTED EXPLORATION ALGORITHM

209 We propose an innovative action-selected exploration for deep RL. Suppose at a predefined time step 210  $t = T_{vs}$ , a set of actions  $\widehat{G}$  is identified from the buffered data, where the cardinality of  $\widehat{G}(|\widehat{G}|)$  is 211 d, with  $d \le p$  indicting a possibly reduced dimension. A critical challenge arises in leveraging the 212 insights gained from action selection for updating the deep RL models. The conventional approach 213 of constructing an entirely new model based on the selected actions is not only time-consuming but also inefficient, particularly in dynamic, non-stationary environments where the requisite action sets 214 are subject to frequent changes. Although our study primarily focuses on stationary environments, 215 the inefficiency of model reinitialization post-selection remains a notable concern.

241 242

254

Algorithm 1 Action-Selected Exploration in Reinforcement Learning
<b>Require:</b> FDR rate $\alpha$ , majority voting ratio $\Gamma$ , max steps T, variable selection step $T_{vs}$
<b>Begin:</b> Initialize the selection set $\widehat{G} = \{\}$ , policy $\pi_{\theta}$ , value function parameter $\phi$ , augmented
replay buffer $\mathcal{D}$
while steps smaller than $T$ do
Sample $\mathbf{a}_{t} \sim \pi_{\theta} \left( \cdot \mid \mathbf{s}_{t} \right)$
Sample knockoff copy $\widetilde{\mathbf{a}}_t \sim \pi_{ heta} \left( \cdot \mid \mathbf{s}_t  ight)$
$\mathbf{s}_{t+1} \sim \mathrm{Env}\left(\mathbf{a}_t, \mathbf{s}_t ight)$
$\mathcal{D} \leftarrow \mathcal{D} \cup \{\mathbf{s}_t, \mathbf{a}_t, \widetilde{\mathbf{a}}_t, r_t, \mathbf{s}_{t+1}\}$
if $t = T_{vs}$ then
Utilize a variable selection algorithm (optional: Knockoff-Sampling in Algorithm 2) on $\mathcal{D}$ to
obtain the estimated minimal sufficient action set $\widehat{G}$
Generate a mask m based on $\widehat{G}$ to prune RL networks based on equation 3 and equation 4
end if
if it's time to update then
update $\phi$ and $\theta$ based on the specific RL algorithm used
end if
end while

235 To seamlessly and efficiently integrate action selection results into deep RL, we propose to mask the 236 non-selected actions and remove their influence once a hard mask is constructed, and thus is flexible 237 to integrate with arbitrary variable selection method. Specifically, in continuous control tasks, deep RL algorithms utilize a policy network  $\pi_{\theta}$  to sample a certain action a given current state s, namely 238  $\mathbf{a} \sim \pi_{\theta} (\cdot | \mathbf{s})$ . The policy network is usually parameterized by a multivariate Gaussian with the 239 diagonal covariance matrix as: 240

 $\mathbf{a} \sim \mathcal{N}\left(\mu\left(\mathbf{s}\right), \operatorname{diag}\left(\sigma\left(\mathbf{s}\right)\right)^{2}\right),$ 

where  $\mu$  and  $\sigma$  are parameterized functions to output mean and standard deviations. Each time we 243 obtain an action from the policy network. The updates of the policy network  $\pi_{\theta}$  and the action-244 value function  $Q_{\phi}$  usually involve sampled actions  $\mathbf{a}_t$  and  $\log \pi_{\theta} (\mathbf{a}_t | \mathbf{s}_t)$  which is the log density 245 of sampled actions. Our strategy is using a binary mask to set them as a certain constant value 246 during the forward pass and it will block the gradient when doing backpropagation and also remove 247 influence when fitting a function. Given a selected action set G, we focus on integrating this selection 248 into the model components  $Q_{\phi}(\mathbf{a}, \mathbf{s})$  and  $\pi_{\theta}(\mathbf{a} \mid \mathbf{s})$ . To facilitate this, we define a selection vector 249  $\mathbf{m} = (m_1, \cdots, m_p) \in \{0, 1\}^p$ , where  $m_i = 1$  if  $i \in \hat{G}$  and 0 otherwise. This vector enables the 250 application of a selection mask to both the  $Q_{\phi}$  and  $\pi_{\theta}$  components as follows. 251

For  $Q_{\phi}$ , we use the hard mask to remove the influence of non-selected actions by not using it in the 252 Q function fitting, 253

$$Q_{\phi}^{m}(\mathbf{a}, \mathbf{s}) = Q_{\phi}(\mathbf{m} \odot \mathbf{a}, \mathbf{s}), \tag{3}$$

where  $\odot$  is the element-wise product. The adoption of action selection significantly reduces the 255 dimensionality of the input action space, thereby reducing bias in the Q function fitting. 256

For  $\pi_{\theta}$ , considering the necessity of updating the policy network via policy gradient, we integrate a 257 hard mask into the logarithm of the policy probability. The modified log probability is formulated as 258 259 lo

$$g \pi_{\theta}^{m}(\mathbf{a} \mid \mathbf{s}) = \mathbf{m} \cdot \left( \log \pi_{\theta}(a_{1} \mid \mathbf{s}), \dots, \log \pi_{\theta}(a_{p} \mid \mathbf{s}) \right), \tag{4}$$

260 where  $\cdot$  is the dot product. This masking of the log probability helps mitigate the likelihood of 261 encountering extremely high entropy values, thereby facilitating a more stable and efficient training 262 process. Hence, this leads to improved model performance. We demonstrate the integration of action 263 selection into deep RL as detailed in Algorithm 1.

264 **Remark 3.1.** Here, we focus on the case that actions are parameterized as diagonal Gaussian which 265 are conditionally independent given states. However, our method can be easily extended to the correlated actions, with details provided in Appendix D. 266

**Remark 3.2.** In scenarios where the algorithm exclusively employs the state-value function  $V_{\phi}(s)$ , 267 the use of the mask operation is unnecessary. Our empirical studies suggest that, even without 268 masking, the model maintains robust performance. This implies that updates to the policy network 269 may hold greater significance than those to the critic in certain contexts.



Figure 2: Learning curves in the Ant-v4 environment reveal that the Knockoff Sampling (KS) method outperforms the traditional Variable Selection (VS) method. When implemented with either the Proximal Policy Optimization (PPO) or Soft-Actor-Critic (SAC) algorithm, KS achieves performance comparable to that of the true actions by automatically setting optimal thresholds to filter out redundant actions, whereas VS often selects useless actions, leading to a high false discovery rate.

-200

0.4 0.6 Total Enviroment Steps (1e6) 0.8

1.0

### 3.2 KNOCKOFF-SAMPLING FOR ACTION SELECTION

0.4 Total Enviro

295

296

297

298

299

300 301 302

303

0.6 ent Steps (1e6) 0.8

Despite the large volume of variable selection (VS) methods [see e.g., 51; 13; 8; 3; 28; 26], these works either *suffer from model-based constraints or lack theoretical guarantees*. The traditional VS often identifies unimportant actions, leading to a high false discovery rate and further causing performance degeneration, as shown in Fig. 2. To provide a general action selection approach for deep RL with false discovery rate control, we propose a novel knockoff-sampling (KS) method that handles dependent data in the online setting with a model-agnostic nature as follows.

310 Suppose now we have a data buffer with the size M, collected from N trajectories where each trajectory has length  $T_j$  for j = 1, ..., N and  $\sum_{j=1}^{N} T_j = M$ . Each time we obtain an action from the 311 312 policy network, we also resample a knockoff copy conditional on the same state  $\tilde{\mathbf{a}}_t \sim \pi_{\theta}$  ( $\cdot | \mathbf{s}_t$ ), and 313 append it to the buffer. The transition tuples thus is redefined as  $(\mathbf{S}_t, \mathbf{A}_t, \mathbf{A}_t, \mathbf{R}_t, \mathbf{S}_{t+1})$ . Note that 314 steps within each trajectory are temporally dependent. To address the issues of highly dependent data 315 in online RL, we couple our method with sample splitting and majority vote following Ma et al. [34]. 316 The proposed KS method consists of three steps as summarized in Algorithm 2: 1. Sample Splitting; 317 2. Knockoff-Sampling Variable Selection; 3. Majority Vote. We detail each step below. 318

1. Sample Splitting: We first split all transition tuples  $(\mathbf{S}_t, \mathbf{A}_t, \mathbf{A}_t, R_t, \mathbf{S}_{t+1})$  into K non-overlapping sub-datasets. This process results in a segmentation of the dataset  $\mathcal{D}$  into distinct subsets  $\mathcal{D}_k$  for  $k \in [K]$ . We combine response variables and denote  $\mathbf{Y}_t = (R_t, \mathbf{S}_{t+1})$  to simplify the notation, based on the target outcomes in Definition 2.1. Here, each sequence  $(\mathbf{S}_t, \mathbf{A}_t, \mathbf{A}_t, \mathbf{Y}_t)$  is assigned to  $\mathcal{D}_k$  if tmod K = k - 1. Subsequent to this division, any two sequences located within the same subset  $\mathcal{D}_k$ either originate from the same trajectory with a temporal separation of no less than K or stem from different trajectories. If the system adheres to  $\beta$ -mixing conditions [6], then a careful selection [5] can allow us to assert that transition sequences within each subset  $\mathcal{D}_k$  is approximately independent.

2. Knockoff-Sampling Variable Selection: For each data subset  $\mathcal{D}_k$ , we select a minimal sufficient 327 action set using the model-X knockoffs as the base selector. Unlike the knockoff method detailed 328 in Section 2.3 that either constructs knockoff features based on second-order machines or estimates the full distribution, we directly sample a knockoff copy of actions from the policy network, i.e., 330  $\widetilde{\mathbf{a}}_t \sim \pi_{\theta}$  ( $\cdot | \mathbf{s}_t$ ). This helps us to bypass the common challenge of creating accurate knockoff features 331 in model-free knockoffs. We theoretically validate that the sampled knockoffs in online RL meet 332 the swapping property equation 1 in Section 4. For every single dimension i of the outcome vector 333  $\mathbf{Y}_t = (R_t, \mathbf{S}_{t+1})$ , we use a general machine learning method (e.g., LASSO, random forest, neural 334 networks) to provide variable importance scores  $Z_{j,i}$  and  $Z_{j,i}$  for the j-th dimension of actions and its knockoff copy, respectively. By the maximum score  $Z_j = \max_i Z_{j,i}$  and  $\widetilde{Z}_j = \max_i \widetilde{Z}_{j,i}$ , the 335 336 selected action set  $\hat{G}_k$  is then obtained following the same procedure in Section 2.3. 337

3. **Majority Vote:** To combine the results on the whole K folds, we calculate the frequency of subsets where the *j*-th action is chosen, i.e.,  $\hat{p}_j = \sum_{k=1}^{K} \mathbb{I}(j \in \hat{G}_k)/K$ , and establish the ultimate selection of actions  $\hat{G} = \{j : \hat{p}_j \ge \Gamma\}$ , with  $\Gamma$  being a predetermined cutoff between 0 and 1.

### 4 THEORETICAL RESULTS

338

339

340 341 342

343 344

345

346

347

348

357

Without loss of generality, we assume that the data buffer  $\mathcal{D}$  consisting of N i.i.d. finite-horizon trajectories, each of length T, which can be summarized as NT transition tuples. We first define two properties to establish theoretical results.

**Definition 4.1.** (Flip Sign Property for Augmented Data) property on the augmented data matrix  $\mathcal{D}_k = \left[\mathbf{A}_k, \tilde{\mathbf{A}}_k, \mathbf{S}_k, \mathbf{Y}_k\right]$  if for any  $j \in [p]$  and  $\Omega \subset [p]$ ,

$$W_i\left(\left[\mathbf{A}_k, \tilde{\mathbf{A}}_k\right]_{\mathrm{swap}(\Omega)}, \mathbf{S}_k, \mathbf{Y}_k\right) = \begin{cases} -W_i\left(\left[\mathbf{A}_k, \tilde{\mathbf{A}}_k\right], \mathbf{S}_k, \mathbf{Y}_k\right), & \text{if } j \in \Omega, \\ W_i\left(\left[\mathbf{A}_k, \tilde{\mathbf{A}}_k\right], \mathbf{S}_k, \mathbf{Y}_k\right), & \text{otherwise}, \end{cases}$$

where  $\mathbf{A}_k, \tilde{\mathbf{A}}_k \in \mathbb{R}^{(NT/K) \times p}$  denote the matrices of the actions and their knockoffs,  $\mathbf{S}_k \in \mathbb{R}^{(NT/K) \times d}$  denote the matrice of state,  $\mathbf{Y}_k \in \mathbb{R}^{(NT/K) \times d+1}$  denotes the response matrix, and  $\begin{bmatrix} \mathbf{A}_k, \tilde{\mathbf{A}}_k \end{bmatrix}_{swap(\Omega)}$  is obtained by swapping all *j*-th columns in  $\mathbf{A}_k, \tilde{\mathbf{A}}_k$  for  $j \in \Omega$ .

The above flip sign property is a common property that needs to be satisfied in knockoff-type methods. We show our method automatically satisfies this property in Lemma F.3 of Appendix.

**Definition 4.2.** (Stationarity and Exponential  $\beta$ -Mixing) The process  $\{(\mathbf{S}_t, \mathbf{A}_t, R_t)\}_{t\geq 0}$  is stationary and exponentially  $\beta$ -mixing if its  $\beta$ -mixing coefficient at time lag k is of the order  $\rho^k$  for some  $0 < \rho < 1$ .

This exponential  $\beta$ -mixing condition has been assumed in the RL literature [see e.g., 2; 10] to derive the theoretical results for the dependent data. Such a condition quantifies the decay in dependence as the future moves farther from the past to achieve the dependence of the future on the past. Based on the above definitions, we establish the following false discovery control results of our method.

Theorem 4.3. Set the number of sample splits  $K = k_0 \log(NT)$  for some  $k_0 > -\log^{-1} \rho$  where  $\rho$  is defined in Definition 4.2. Assume that the following assumption hold: the process  $\{(\mathbf{S}_t, \mathbf{A}_t, R_t)\}_{t \ge 0}$ is stationary and exponentially  $\beta$ -mixing.

Then  $\widehat{G}_k$  obtained by Algorithm 2 with the standard knockoffs controls the modified FDR (mFDR), mFDR  $\leq \alpha + O\{K^{-1}(NT)^{-c}\},\$ 

373 where the constant  $c = -k_0 \log(\rho) - 1 > 0$ .

The proof can be mainly divided into two parts, firstly we show valid mFDR control can be achieved when data are independent, then for dependent data satisfying the  $\beta$ -mixing condition, the upper band can be relaxed as the cost of dependence. Finally, a combination of the two would provide the final upper bound on mFDR control. The detailed proof is in Appendix F.

# <sup>378</sup> 5 EXPERIMENTS

379

380 **Experiment Setup** We aim to answer whether the variable selection is helpful for deep RL training when the action dimension is high and redundant. We conduct experiments on standard locomotion 382 tasks in MuJoCo [52] and treatment allocation tasks calibrated from electronic health records (EHR), the MIMIC-III dataset [19]. The environment details are in Table B.1 of Appendix. Here, we focus on two representative actor-critic algorithms, Proximal Policy Optimization (PPO) [43] and 384 Soft-Actor-Critic (SAC) [15]. We adopt the implementation from Open AI Spinning up Framework 385 [1]. For SAC, the implementation involves fitting both  $Q_{\phi}$  and  $\pi_{\theta}$ , and we use a mask on both 386 components. For PPO, it fits  $V_{\theta}$  and  $\pi_{\theta}$ , hence we only combine the mask with  $\pi_{\theta}$ . Tables A.1 and 387 A.2 summarize the hyperparameters we used. We set the FDR rate  $\alpha = 0.1$  and voting ratio  $\Gamma = 0.5$ 388 in all settings. All the experiments are conducted in the server with  $4 \times$  NVIDIA RTX A6000 GPU. 389

Semi-synthetic MuJoCo Environments We choose three tasks: Ant, HalfCheetah, and Hopper. To 390 increase the dimension of action space, we artificially add extra p actions to the raw action space and 391 consider two scenarios p = 20 and 50. For each setting, we run experiments over  $2 \times 10^5$  and  $10^6$ 392 steps for SAC and PPO, respectively, averaged over 10 training runs. The running steps for SAC and 393 PPO are set adaptively to obtain better exploration for each method and save computation costs, as the 394 main goal is to show how action selection can improve sample efficiency rather than compare these 395 two methods. For each evaluation point, we run 10 test trajectories and average their reward as the average return. Besides RL algorithm performance, we also evaluate variable selection performance 397 in terms of True Positive Rate (TPR), False Positive Rate (FPR), and FDR. 398

Action Selection in the Initial Stage of Training We utilize action selection in the beginning 399 stage of the training. For both methods, we utilize the first 4000 samples for variable selection 400 and then use the selection results to build a hard mask for action in deep RL models. We compare 401 our knockoff sampling (KS) method with the baseline of selecting all actions (All) to evaluate the 402 impact of integrating a masking mechanism with a selection strategy in deep RL. We also provide 403 the experimental results with only ground-truth actions selected (True) as a reference. To reduce 404 the computational complexity, we choose LASSO [51] as our base variable selection algorithm 405 for KS. Here, selecting all actions (All) and ground-truth actions (True) are the cases where RL 406 models trained on whole action space and minimal sufficient action space, respectively. Hence, the 407 model corresponding to ground-truth action has smaller parameters than all other methods because its initialization is based on the minimal sufficient action set. The results are shown in Fig. 3 for 408 PPO, Fig. B.1 for SAC, and Table 1 for all numerical details. Due to space constraints, we mainly 409 present the PPO figures in the main text. In all cases, we find that KS-guided models outperform 410 those trained on the whole action space in terms of average return and much lower FDR and FPR, 411 with larger improvement gains as p increases. This empirically validates our theory of FDR control 412 with the proposed KS method, demonstrating that action selection can enhance learning efficiency 413 during the initial stages of RL training where action space is high and redundant. 414

Action Selection in the Middle Stage of Training To show whether action selection can be used 415 in the middle stage of training to remedy the inefficiency brought by exploring the whole action 416 space, we conduct experiments where in the first half of the training steps the models are trained 417 on the whole action space, and in the middle of the stage, we utilize action selection and build hard 418 masks for them and then continue training for the rest of the steps. We compare our KS method 419 with selecting all actions (All) and ground-truth actions (True) similarly. The results in Fig. 4 and 420 Fig. B.2 reveal a notable pattern: agents initially struggle to learn effectively, but mid-stage variable 421 selection significantly improves their performance, with models trained on the correct actions. This 422 demonstrates the effectiveness of mid-stage variable selection in enhancing learning outcomes.

423 Treatment Allocation for Sepsis Patients We test our method with PPO and utilize the first 1000 424 samples for action selection in the initial stage. We also include two additional baselines: Lattice 425 [9], and gSDE [38]. Lattice and gSDE also use all actions but additionally incorporate temporally 426 correlated Gaussian noise into the training. We run experiments over  $5 \times 10^4$  time steps and average 427 results over 5 runs. For each evaluation point, we run 5 test trajectories and average their reward 428 as the average return. The results are shown in Table 2 and Fig. 5. Lattice and gSDE exhibit high 429 variance and instability in this environment, likely due to over-exploration. In contrast, our method consistently delivers stable and superior performance compared to the others. The analysis in Fig. 5 430 reveals that our methods effectively identify sepsis-influencing treatments that are relevant to key 431 body factors, demonstrating the potential for real-world medical applications.



Figure 3: Learning curves for PPO in the MujoCo environments with different approaches during the initial stage. In all experiments, our knockoff sampling (KS) method not only performs comparably to the true actions but also consistently delivers higher rewards than using all actions.



Figure 4: Learning curves for PPO in the MujoCo environments during the middle stage where the red line indicates the time point we utilize the proposed KS. After identifying the essential action set, the policy can be more efficient and achieve higher rewards than continuing training on all actions.

Action Selection is Fast and Lightweight With just a few thousand data points and a lightweight machine learning algorithm like random forest or LASSO, the whole action selection process outlined in Algorithm 2, completes in under 20 seconds—including knockoff threshold determination. This is significantly faster and less computationally intensive than the RL training part. Even when incorporating more sophisticated feature selection methods, the additional computational overhead remains *negligible* compared to the time required for RL training. Moreover, for the RL agent's deep neural network, only a few lightweight masking parameters are introduced, which have minimal effect on both training and inference speed. Yet, these in turn substantially enhance policy optimization.

**Supporting Analyses** We conduct additional experiments for the variation of action distributions during training, the convergence behavior of PPO with more steps, and whether network capability is a key factor in solving the high-dimension action problem. The details are in Appendix C.

## 6 CONCLUSION, LIMITATION, AND FUTURE WORK

In this work, we address the high-dimensional action selection problem in online RL. We formally define the objective of action selection by identifying a minimal sufficient action set. We innovate by integrating a knockoff-sampling variable selection into broadly applicable deep RL algorithms. Empirical evaluations in synthetic robotics and treatment allocation environments demonstrate the enhanced efficacy of our approach. Yet, a notable constraint of our method is its singular application during the training phase, coupled with the potential risk of overlooking essential actions with weak signals. Inadequate action selection could degrade the agent's performance. Intriguing future research includes extending our methodology to incorporate multiple and adaptive selection stages. This adaptation could counterbalance initial omissions in action selection. Additionally, formulating an effective termination criterion for this process represents another compelling research direction.

Table 1: Results on the PPO and SAC for three Mujoco tasks: Ant, HalfCheetah, and Hopper. Action selection is utilized at the beginning stage of RL training. The final reward is the performance evaluation for the agent after training. The best-performing results between KS and All are highlighted in bold.

Env	RL Algo	n	Selection	Ant			
	itte ringo.	Р	Selection	TPR $(\uparrow)$	FDR $(\downarrow)$	FPR $(\downarrow)$	Reward (†)
		0	True	1.00	0.0	0.00	567.77
			KS	1.00	0.01	0.01	507.90
	PPO	20	All	1.00	0.71	1.00	202.65
			KS	1.00	0.00	0.00	572.39
Δnt		50	All	1.00	0.86	1.00	151.66
2 111		0	True	1.00	0.00	0.00	817.95
			KS	1.00	0.01	0.01	937.74
	SAC	20	All	1.00	0.71	1.00	12.61
			KS	1.00	0.00	0.00	731.73
		50	All	1.00	0.86	1.00	-208.04
		0	True	1.00	0.0	0.00	2130.55
			KS	1.00	0.01	0.01	2237.08
	PPO	20	All	1.00	0.77	1.00	1356.46
			KS	1.00	0.00	0.00	1932.27
HalfCheetah		50	All	1.00	0.89	1.00	619.67
Hancheetan		0	True	1.00	0.00	0.00	6640.05
			KS	1.00	0.00	0.00	6607.55
	SAC	20	All	1.00	0.77	1.00	5631.20
			KS	1.00	0.00	0.00	6873.95
		50	All	1.00	0.89	1.00	4748.24
	РРО	0	True	1.00	0.0	0.00	1736.65
			KS	1.00	0.00	0.00	1540.83
		20	All	1.00	0.87	1.00	1205.12
			KS	1.00	0.00	0.00	1710.82
Honnar		50	All	1.00	0.94	1.00	703.08
riopper		0	True	1.00	0.00	0.00	2511.00
	SAC		KS	1.00	0.00	0.00	2165.81
		20	All	1.00	0.87	1.00	398.75
			KS	1.00	0.00	0.00	2424.09
		50	All	1.00	0.94	1.00	137.67

Table 2: In the treatment allocation environments, we report the average reward and standard deviation (Std) at both the midpoint (50%) and the end (100%) of the training process, where the proposed KS performs the best with the highest reward. Both KS and All show significantly lower variance with improving performance over time. In contrast, Lattice and gSDE show high variance and degraded performance.

Method	Average Reward (Std)		
wichiou	50%	100%	
KS	14.0(0.8)	<b>15</b> .3(1.2)	
All	13.6(1.6)	14.4(1.0)	
Lattice	12.7(1.7)	10.8(5.9)	
gSDE	9.7(5.3)	9.3(5.3)	

Actions (after states) Treatment Index

Figure 5: In the treatment allocation environments, our KS methods identify the most relevant treatments and the policy can be more efficient to achieve higher rewards than continuing training on all actions.

540	REFERENCES
541	THE BILLICED

544

546

547

549

550

551

552 553

554

555

556

558 559

561

562

563 564

565

566

567

568

569

570

571

572 573

574

575

576

577 578

579

580

581

582

583

584 585

586

587

588

589

590

591

- [1] Achiam, J. Spinning Up in Deep Reinforcement Learning. 2018.
- [2] Antos, A., Szepesvári, C., and Munos, R. Learning near-optimal policies with bellman-residual minimization based fitted policy iteration and a single sample path. Machine Learning, 71: 89-129, 2008.
- [3] Balın, M. F., Abid, A., and Zou, J. Concrete autoencoders: Differentiable feature selection and reconstruction. In International conference on machine learning, pp. 444–453. PMLR, 2019. 548
  - [4] Barber, R. F. and Candès, E. J. Controlling the false discovery rate via knockoffs. 2015.
    - [5] Berbee, H. Convergence rates in the strong law for bounded mixing sequences. Probability theory and related fields, 74(2):255–270, 1987.
      - [6] Bradley, R. C. Basic properties of strong mixing conditions. a survey and some open questions. 2005.
    - [7] Candes, E., Fan, Y., Janson, L., and Lv, J. Panning for gold: 'model-x'knockoffs for high dimensional controlled variable selection. Journal of the Royal Statistical Society Series B: Statistical Methodology, 80(3):551–577, 2018.
    - [8] Chen, J., Stern, M., Wainwright, M. J., and Jordan, M. I. Kernel feature selection via conditional covariance minimization. Advances in Neural Information Processing Systems, 30, 2017.
    - [9] Chiappa, A. S., Marin Vargas, A., Huang, A., and Mathis, A. Latent exploration for reinforcement learning. Advances in Neural Information Processing Systems, 36, 2024.
    - [10] Dai, B., Shaw, A., Li, L., Xiao, L., He, N., Liu, Z., Chen, J., and Song, L. Sbeed: Convergent reinforcement learning with nonlinear function approximation. In International conference on machine learning, pp. 1125–1134. PMLR, 2018.
    - [11] Farquhar, G., Gustafson, L., Lin, Z., Whiteson, S., Usunier, N., and Synnaeve, G. Growing action spaces. In International Conference on Machine Learning, pp. 3040–3051. PMLR, 2020.
    - [12] Gershman, S. J., Pesaran, B., and Daw, N. D. Human reinforcement learning subdivides structured action spaces by learning effector-specific values. Journal of Neuroscience, 29(43): 13524–13531, 2009.
    - [13] Gu, Q., Li, Z., and Han, J. Generalized fisher score for feature selection. arXiv preprint arXiv:1202.3725, 2012.
    - [14] Guo, Z. D. and Brunskill, E. Sample efficient feature selection for factored mdps. arXiv preprint arXiv:1703.03454, 2017.
    - [15] Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In International conference on machine learning, pp. 1861–1870. PMLR, 2018.
    - [16] Hao, B., Duan, Y., Lattimore, T., Szepesvári, C., and Wang, M. Sparse feature selection makes batch reinforcement learning more sample efficient. In International Conference on Machine Learning, pp. 4063-4073. PMLR, 2021.
    - [17] Jain, V., Fedus, W., Larochelle, H., Precup, D., and Bellemare, M. G. Algorithmic improvements for deep reinforcement learning applied to interactive fiction. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pp. 4328–4336, 2020.
    - [18] Jiang, Y., Van Wouwe, T., De Groote, F., and Liu, C. K. Synthesis of biologically realistic human motion using joint torque actuation. ACM Transactions On Graphics (TOG), 38(4): 1-12, 2019.
  - [19] Johnson, A. E., Pollard, T. J., Shen, L., Lehman, L.-w. H., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Anthony Celi, L., and Mark, R. G. Mimic-iii, a freely accessible critical care database. Scientific data, 3(1):1-9, 2016.

600 601

602

603

604

605 606

607

608

609

610

611

612 613

614

615

616

617

618

619

620 621

622

623

624

625 626

627

628

629

630

631

632 633

634

635

636 637

638

639 640

641

642 643

644

- [20] Jordon, J., Yoon, J., and van der Schaar, M. Knockoffgan: Generating knockoffs for feature selection using generative adversarial networks. In International conference on learning representations, 2018.
  - [21] Kaiser, L., Babaeizadeh, M., Milos, P., Osinski, B., Campbell, R. H., Czechowski, K., Erhan, D., Finn, C., Kozakowski, P., Levine, S., et al. Model-based reinforcement learning for atari. arXiv preprint arXiv:1903.00374, 2019.
  - [22] Kanervisto, A., Scheller, C., and Hautamäki, V. Action space shaping in deep reinforcement learning. In 2020 IEEE conference on games (CoG), pp. 479–486. IEEE, 2020.
  - [23] Kober, J., Bagnell, J. A., and Peters, J. Reinforcement learning in robotics: A survey. <u>The</u> International Journal of Robotics Research, 32(11):1238–1274, 2013.
  - [24] Kolm, P. N. and Ritter, G. Modern perspectives on reinforcement learning in finance. <u>Modern Perspectives on Reinforcement Learning in Finance (September 6, 2019)</u>. The Journal of Machine Learning in Finance, 1(1), 2020.
  - [25] Kroon, M. and Whiteson, S. Automatic feature selection for model-based reinforcement learning in factored mdps. In <u>2009 International Conference on Machine Learning and Applications</u>, pp. 324–330. IEEE, 2009.
  - [26] Lee, C., Imrie, F., and van der Schaar, M. Self-supervision enhanced feature selection with correlated gates. In International Conference on Learning Representations, 2021.
  - [27] Lee, K., Kim, S.-A., Choi, J., and Lee, S.-W. Deep reinforcement learning in continuous action spaces: a case study in the game of simulated curling. In <u>International conference on machine</u> learning, pp. 2937–2946. PMLR, 2018.
  - [28] Liang, F., Li, Q., and Zhou, L. Bayesian neural networks for selection of drug sensitive genes. Journal of the American Statistical Association, 113(523):955–972, 2018.
  - [29] Liu, J., Sun, A., and Ke, Y. A generalized knockoff procedure for fdr control in structural change detection. Journal of Econometrics, 2022.
  - [30] Liu, Y. and Zheng, C. Auto-encoding knockoff generator for fdr controlled variable selection. arXiv preprint arXiv:1809.10765, 2018.
  - [31] Liu, Y., Logan, B., Liu, N., Xu, Z., Tang, J., and Wang, Y. Deep reinforcement learning for dynamic treatment regimes on medical registry data. In <u>2017 IEEE international conference on</u> healthcare informatics (ICHI), pp. 380–385. IEEE, 2017.
  - [32] Luo, S., Androwis, G., Adamovich, S., Nunez, E., Su, H., and Zhou, X. Robust walking control of a lower limb rehabilitation exoskeleton coupled with a musculoskeletal model via deep reinforcement learning. Journal of neuroengineering and rehabilitation, 20(1):1–19, 2023.
  - [33] Ma, S., Dalgleish, J., Lee, J., Wang, C., Liu, L., Gill, R., Buxbaum, J. D., Chung, W. K., Aschard, H., Silverman, E. K., et al. Powerful gene-based testing by integrating long-range chromatin interactions and knockoff genotypes. <u>Proceedings of the National Academy of Sciences</u>, 118 (47):e2105191118, 2021.
  - [34] Ma, T., Cai, H., Qi, Z., Shi, C., and Laber, E. B. Sequential knockoffs for variable selection in reinforcement learning. arXiv preprint arXiv:2303.14281, 2023.
  - [35] Misra, D., Henaff, M., Krishnamurthy, A., and Langford, J. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. In <u>International conference on</u> machine learning, pp. 6961–6971. PMLR, 2020.
  - [36] Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. Playing atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- [37] Pavse, B. S. and Hanna, J. P. Scaling marginalized importance sampling to high-dimensional state-spaces via state abstraction. In <u>Proceedings of the AAAI Conference on Artificial</u> Intelligence, volume 37, pp. 9417–9425, 2023.

652 653

654

655 656

657

658

659

660

661

662

663

664 665

666

667

668

669

670

671 672

673

674 675

676

677

678

679

680

681 682

683

684 685

686

687

688

689

690

691

692

693 694

695

696 697

698

- [38] Raffin, A., Kober, J., and Stulp, F. Smooth exploration for robotic reinforcement learning. In Conference on Robot Learning, pp. 1634–1644. PMLR, 2022.
  - [39] Romano, Y., Sesia, M., and Candès, E. Deep knockoffs. Journal of the American Statistical Association, 115(532):1861–1872, 2020.
  - [40] Ruan, S., Comanici, G., Panangaden, P., and Precup, D. Representation discovery for mdps using bisimulation metrics. In <u>Proceedings of the AAAI Conference on Artificial Intelligence</u>, volume 29, 2015.
  - [41] Sadamoto, T., Chakrabortty, A., and Imura, J.-i. Fast online reinforcement learning control using state-space dimensionality reduction. <u>IEEE Transactions on Control of Network Systems</u>, 8(1):342–353, 2020.
  - [42] Sakryukin, A., Raïssi, C., and Kankanhalli, M. Inferring dqn structure for high-dimensional continuous control. In <u>International Conference on Machine Learning</u>, pp. 8408–8416. PMLR, 2020.
  - [43] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
  - [44] Schumacher, P., Haeufle, D., Büchler, D., Schmitt, S., and Martius, G. DEP-RL: Embodied exploration for reinforcement learning in overactuated and musculoskeletal systems. In <u>The Eleventh International Conference on Learning Representations</u>, 2023. URL https://openreview.net/forum?id=C-xa\_D3oTj6.
  - [45] Sesia, M., Sabatti, C., and Candès, E. J. Gene hunting with knockoffs for hidden markov models. arXiv preprint arXiv:1706.04677, 2017.
  - [46] Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. nature, 529(7587):484–489, 2016.
  - [47] Sun, X., Mao, T., Ray, L., Shi, D., and Kralik, J. Hierarchical state-abstracted and socially augmented q-learning for reducing complexity in agent-based learning. <u>Journal of Control</u> Theory and Applications, 9:440–450, 2011.
  - [48] Sunehag, P., Evans, R., Dulac-Arnold, G., Zwols, Y., Visentin, D., and Coppin, B. Deep reinforcement learning with attention for slate markov decision processes with high-dimensional states and actions. arXiv preprint arXiv:1512.01124, 2015.
  - [49] Synnaeve, G., Gehring, J., Lin, Z., Haziza, D., Usunier, N., Rothermel, D., Mella, V., Ju, D., Carion, N., Gustafson, L., et al. Growing up together: Structured exploration for large action spaces. 2019.
    - [50] Tan, H., Zhang, H., Peng, J., Jiang, Z., and Wu, Y. Energy management of hybrid electric bus based on deep reinforcement learning in continuous state and action space. <u>Energy Conversion</u> and Management, 195:548–560, 2019.
  - [51] Tibshirani, R. Regression shrinkage and selection via the lasso. Journal of the Royal Statistical Society Series B: Statistical Methodology, 58(1):267–288, 1996.
  - [52] Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. In <u>2012</u> <u>IEEE/RSJ international conference on intelligent robots and systems</u>, pp. 5026–5033. IEEE, 2012.
  - [53] Xiao, B., Lu, Q., Ramasubramanian, B., Clark, A., Bushnell, L., and Poovendran, R. Fresh: Interactive reward shaping in high-dimensional state spaces using human feedback. <u>arXiv</u> preprint arXiv:2001.06781, 2020.
  - [54] Yu, C., Liu, J., Nemati, S., and Yin, G. Reinforcement learning in healthcare: A survey. <u>ACM</u> Computing Surveys (CSUR), 55(1):1–36, 2021.
- [55] Zahavy, T., Haroush, M., Merlis, N., Mankowitz, D. J., and Mannor, S. Learn what not to learn: Action elimination with deep reinforcement learning. <u>Advances in neural information</u> processing systems, 31, 2018.

#### **IMPLEMENTATION DETAILS** А

We adopt the implementation from Open AI Spinning up Framework [1]. Tables A.1 and A.2 show the hyperparameters for the RL algorithms we used in our experiments. We set the FDR rate  $\alpha = 0.1$ and voting ratio r = 0.5 for our knockoff method in all settings. All the experiments are conducted in the server with  $4 \times$  NVIDIA RTX A6000 GPU.

710			
711	Parameter	Mujoco	EHR
712	optimizer	Adam	Adam
713	learning rate $\pi$	$3.0 \cdot 10^{-4}$	$3.0 \cdot 10^{-3}$
714	learning rate V	$1.0 \cdot 10^{-3}$	$1.0 \cdot 10^{-3}$
715	learning rate schedule	constant	constant
716	discount ( $\gamma$ )	0.99	0.99
717	number of hidden layers (all networks)	2	2
718	number of hidden units per layer	[64, 32]	[64, 32]
710	number of samples per minibatch	256	100
719	number of steps per rollout	1000	100
720	non-linearity	ReLU	ReLU
721	gSDE		
722	initial log $\sigma$	0	0
723	Full std matrix	Yes	Yes
724	Lattice		
725	initial log $\sigma$	0	0
726	Full std matrix	Yes	Yes
727	Std clip	(0.001,1)	(0.001,1)

### Table A.1: PPO Hyperparameters

### Table A.2: SAC Hyperparameters

Parameter	Mujoco
optimizer	Adam
learning rate $\pi$	$3.0\cdot10^{-4}$
learning rate Q	$3.0 \cdot 10^{-4}$
learning rate schedule	constant
discount $(\gamma)$	0.9
replay buffer size	$1 \cdot 10^6$
number of hidden layers (all networks)	2
number of hidden units per layer	[256, 256]
number of samples per minibatch	256
non-linearity	ReLU
entropy coefficient ( $\alpha$ )	0.2
warm-up steps	$1.0\cdot 10^4$

### 

#### MORE EXPERIMENTAL RESULTS AND ANALYSES В

### Table B.1: Summary of Environments

Env	Dimension of Action	Dimension of State
Ant	8	27
HalfCheetah	6	17
Hopper	3	11
EHR	20	46



### We list the dimension of action and state in terms of environments we used in Table B.1.

Figure B.1: Results of SAC and PPO when using different variable selection approaches during the initial stage.

### B.1 MuJoCo

The results for the initial stage are shown in Fig. B.1 and Table 1. For different environments, action selection difficulty varies and Hopper is the easiest one where the proposed KS method can correctly select the minimal sufficient action set. Also, the patterns of the results for PPO and SAC are similar. One observation is that KS can select all the sufficient actions with all TPRs equal to 1. It is shown that KS is efficient in selecting only the minimal sufficient action set in almost all scenarios, which also empirically validates our theory of FDR control under the proposed method.

**B.2** TREATMENT ALLOCATION FOR SEPSIS PATIENTS

We collect 250000 data points from the MIMIC-III Clinical Database. Then we utilize a long short-term memory (LSTM) to model the state transition. The observed state information encompasses a
broad spectrum of clinical and laboratory variables for assessing patient health and outcomes in a
medical setting. It includes demographic information (gender, age), physiological metrics (weight, vital signs such as heart rate, blood pressure, respiratory rate, oxygen saturation, temperature), and neurological status (Glasgow Coma Scale). Additionally, it captures details about the patient's



Figure B.2: Learning curves in the MujoCo environments with our method during the middle stage where the red line indicates the time we utilize KS. After identifying a less redundant action set, the policy can be more efficient and achieve higher rewards than continuing training on all actions.

readmission status, mechanical ventilation use, and severity scores such as SOFA and SIRS. Comprehensive lab tests cover a variety of blood chemistry components, including electrolytes, liver enzymes, and arterial blood gases.

The action includes treatments such as the median and maximum doses of vasopressors administered to manage blood pressure and perfusion, alongside vasopressors and intravenous fluids. The reward is calculated based on whether the patient's SOFA Score has improved. Termination of an episode is achieved based on the patient's mortality rate reaching the minimum (SOFA Score being 0) or the patient's mortality rate reaching the maximum (SOFA score being 24).

C SUPPORTING ANALYSES

We conducted additional experiments regarding the distribution of actions that were sampled over the
 training period. The new results are summarized in Figure C.1. Based on the results together with
 existing figures in the main text, we can conclude that there exist changes regarding the distribution
 of actions that are sampled during different periods of training which could be a positive indicator of
 a more focused and potentially more effective learning process.



Figure C.1: The distributions of actions in 3 stages: initial training, middle of training without KS, and middle of training with KS. It can be seen that with KS the actions have slightly less variance than other methods, which could be a positive indicator of a more focused and potentially more effective learning process.







Figure C.3: Learning curves in the MujoCo environments with different network sizes.



918 existing figures in the main text, we can conclude that the benefit of the proposed framework is both 919 sample efficiency and performance. 920

We also conducted additional experiments by increasing the network size. The new results under 921 MujoCo with different network sizes are summarized in Figure C.3, where we can conclude that 922 the network capacity has a certain influence on the performance of All Action. However, simply 923 increasing the network capacity unnecessarily makes learning easier. In addition, the proposed 924 method consistently performs better than All Action no matter how we increase the network size, 925 which indicates the performance difference results from the redundancy of action space. In addition, 926 we admit there are other factors that affect how the agent learns with a larger action dimension, such 927 as regularization techniques, albeit with limited influence compared with the redundancy of action 928 space.

931 932

933

934

935

936

937

938

939

941

946

947

948 949

950

951 952

953

954

955

956

957

958 959

960 961

967

#### EXTEND TO CORRELATED ACTIONS D

Now assume the policy network is parameterized by a multivariate Gaussian with covariance matrix as:

$$\mathbf{a}\sim\mathcal{N}\left(\boldsymbol{\mu}\left(\mathbf{s}\right),\boldsymbol{\Sigma}\left(\mathbf{s}\right)\right),$$

where  $\mu$  and  $\Sigma$  are parameterized functions to output the mean and covariance matrix.

Since the actions are correlated, we couldn't mask the individual  $\log \pi_{\theta} (a_i \mid \mathbf{s})$ . To solve this problem, we can transform the non-selected actions to be conditional independent first.

Given a selected action set  $\widehat{G}$ , we define a selection vector  $\mathbf{m} = (m_1, \dots, m_p) \in \{0, 1\}^p$ , where 940  $m_i = 1$  if  $i \in \hat{G}$  and 0 otherwise. Then we mask the covariance matrix as follows:

$$\Sigma^m(\mathbf{s})_{ij} \begin{cases} \Sigma(\mathbf{s})_{ij} & \text{if } m_i = 1, \text{ and } m_j = 1, i \neq j \\ 0 & \text{if } m_i = 0, \text{ or } m_j = 0, i \neq j. \end{cases}$$

The masking only changes non-selected actions to be independent and removes its influence on selected actions which still keeps the covariance structure of selected actions. Then we can easily mask the log density of non-selected actions.

### MACHINE LEARNING ALGORITHM FOR CALCULATING IMPORTANCE E SCORES

In Knockoff, the computation of the importance scores is very flexible since many machine learning algorithms can be used. The only requirement for the ML method is to satisfy a fairness constraint, so that that swapping  $X_i$  with  $X_j$  would have the only effect of swapping  $Z_j$  with  $Z_j$ , which is usually true in standard tabular machine learning algorithms, like regression, decision tree. One example is that when we fit a linear regression model, then the coefficients of the variables can be seen as importance scores.

#### F TECHNICAL PROOFS

#### PRELIMINARY RESULTS F.1 962

963 Before we prove Theorem 1, we first provide a preliminary lemma of our procedure that can enable the 964 flip-sign property of W-statistics. This property can be used to prove Theorem 1 when observations 965 are independent. Now we focus on one data split  $\mathcal{D}_k$  and assume the data are independent. 966

**Lemma F.1.**  $A_k$  and  $S_k$  are a action and a state matrix. For any subset  $\Omega \subset \{1, \ldots, p\}$ , and  $A_k$ obatined by resampling, we have

$$\left( \left[ \mathbf{A}_k, \tilde{\mathbf{A}}_k \right]_{\mathrm{swap}(\Omega)}, \mathbf{S}_k \right) \stackrel{d}{=} \left( \left[ \mathbf{A}_k, \tilde{\mathbf{A}}_k \right], \mathbf{S}_k \right),$$

where swap( $\Omega$ ) represents swapping the *j*-th entry of  $\mathbf{A}_k$  and  $\mathbf{A}_k$  for all  $j \in \Omega$ .

The proof of Lemma F.1 is based on the property of constructed variables where  $\mathbf{A}_k$  and  $\tilde{\mathbf{A}}$  have the same marginal distribution and the whole joint distribution is symmetrical in terms of  $\mathbf{A}_k$  and  $\tilde{\mathbf{A}}$ .

In the following, we show the exchangeability holds jointly on actions, states, and rewards when
 swapping null variables.

**Lemma F.2.** Let  $\mathcal{H}_0 \subseteq \{1, \ldots, p\}$  be the indices of the null variables, for any subset  $\Omega \subset \mathcal{H}_0$ 

$$\left( \left[ \mathbf{A}_{k}, \tilde{\mathbf{A}}_{k} 
ight]_{\mathrm{swap}(\Omega)}, \mathbf{S}_{k}, \mathbf{Y}_{k} 
ight) \stackrel{d}{=} \left( \left[ \mathbf{A}_{k}, \tilde{\mathbf{A}}_{k} 
ight], \mathbf{S}_{k}, \mathbf{Y}_{k} 
ight),$$

where  $\mathbf{Y}_k$  is a response including the next state and reward.

**Proof:** Based on the exchangeability proved in Lemma F.1, we can directly utilize the proof of Lemma 3.2 in Candès et al. (2018). We just need to extend the derivation by conditional on  $S_k$  and show equivalence by swapping action variables in  $\Omega$  one by one. We omit further details of the proof. **Lemma F.3.** 

Lemma r.s.

$$W_i\left(\left[\mathbf{A}_k, \tilde{\mathbf{A}}_k\right]_{\mathrm{swap}(\Omega)}, \mathbf{S}_k, \mathbf{Y}_k\right) = W_i\left(\left[\mathbf{A}_k, \tilde{\mathbf{A}}_k\right], \mathbf{S}_k, \mathbf{Y}_k\right) \cdot \begin{cases} -1, & \text{if } i \in \Omega \\ +1, & \text{otherwise} \end{cases}$$

**Proof:** We require the method constructing W to satisfy a fairness requirement so that swapping two variables would have the only effect of swapping corresponding feature importance scores. The fairness constraint is satisfied with many general machine learning algorithms, like LASSO and random forest. Once the fairness constraint is satisfied, W will be anti-symmetric and the equal above automatically holds.

**Lemma F.4.** Assume the flip-coin property in Lemma F.3 is satisfied, on data  $\mathcal{D}_k$ , the selection  $\widehat{G}_k$  obtained from applying knockoff method in Algorithm 2 controls modified FDR (mFDR), e.g.

$$mFDR\left(\widehat{G}_k\right) \leq \alpha.$$

**Proof:** For statistics W calculated, we denote  $W_{\text{swap}(\Omega)}$  to be the W-statistics computed after the swap w.r.t.  $\Omega \subset \{1, \ldots, p\}$ . Now consider a sign vector  $\epsilon \in \{\pm 1\}^p$  independent of  $\mathbf{W} = [W_1, \ldots, W_p]^\top$ , where  $\epsilon_i = 1$  for all non-null state variables and  $\mathbb{P}(\epsilon_i = 1) = 1/2$  are independent for all null state variables. Then for such  $\epsilon$ , denote  $\Omega := \{i : \epsilon_i = -1\}$ , which is a subset of  $\mathcal{H}_0$  by the assumption (and recall that  $\mathcal{H}_0$  is the collection of all null variables). By Lemma F.3 we know

$$(W_1 \cdot \epsilon_1, \ldots, W_p \cdot \epsilon_p) = W_{\operatorname{swap}(\Omega)}.$$

For convenience, we also use h to denote a measurable mapping function from a data set to its W-statistics, i.e., on  $[\mathbf{s}_k, \tilde{\mathbf{s}}_k, \mathbf{a}_k, \mathbf{y}_k]$ ,

$$\mathbf{W} = h\left(\mathbf{A}_k, \tilde{\mathbf{A}}_k, \mathbf{S}_k, \mathbf{Y}_k
ight).$$

<sup>1012</sup> Then we can get:

$$egin{aligned} \mathbf{W}_{ ext{swap}(\Omega)} &= h\left(\left[\mathbf{A}_k, ilde{\mathbf{A}}_k
ight]_{ ext{swap}(\Omega)}, \mathbf{S}_k, \mathbf{Y}_k
ight) \ &\stackrel{d}{=} h\left(\left[\mathbf{A}_k, ilde{\mathbf{A}}_k
ight], \mathbf{S}_k, \mathbf{Y}_k
ight) = \mathbf{W}, \end{aligned}$$

where the second equality (in distribution) is due to Lemma F.2 and h is measurable. The rest of the proof will be the same as that for Theorems 1 and 2 in Barber & Candès [4].

1021 F.2 PROOF OF THEOREM 4.3

1023 Using Lemma F.4, we can show that if the data points in  $\mathcal{D}_k$  are independent, then mFDR can be 1024 controlled. Now we want to weaken the independence assumption to stationarity and exponential 1025  $\beta$ -mixing assumption in 4.2. Based on Lemma F.4, the following proof is essentially the same as 1026 Theorem 1 in Ma et al. [34]. We will omit those steps for brevity.