

A Non-Monolithic Policy Approach of Offline-to-Online Reinforcement Learning

JaeYoon Kim¹[0000–0003–0898–5097], Junyu Xuan¹[0000–0002–8367–6908],
Christy Liang²[0000–0001–7179–5208], and Farookh Hussain¹[0000–0003–1513–8072]

¹ Australian Artificial Intelligence Institute (AII), University of Technology Sydney, Australia
JaeYoon.Kim@student.uts.edu.au,
{Junyu.Xuan, Farookh.Hussain}@uts.edu.au
² Visualisation Institute, University of Technology Sydney, Australia
Jie.Liang@uts.edu.au

A Summary of reference paper [1]

The reference paper [1] is summarized in this research.

A.1 Exploration modes

- Exploit mode (G) : the greedy pursuit of external reward
- Explore mode (X) : uniform random (XU), intrinsic reward based on random network distillation (XI)

A.2 Granularity

Four choices of temporal granularity are considered for exploratory periods.

- Step-level : each step
- Experiment-level : all behaviours produced in explore mode during off-policy training
- Episode-level : each episode
- Intra-episodic : between step-level and episode-level

A.3 Switching mechanisms

Four choices of temporal granularity are considered for exploratory periods. ‘Homeostasis’ (See A.5) leverages the ‘value promise discrepancy’, $D_{promise}(t-k, t)$, which represents the difference in the value function over k steps:

$$D_{promise}(t-k, t) := \left| V(s_{t-k}) - \sum_{i=0}^{k-1} \gamma^i R_{t-i} - \gamma^k V(s_t) \right| \quad (1)$$

where $V(s)$ is an agent’s value estimate at state s , R is a reward and γ is a discount factor.

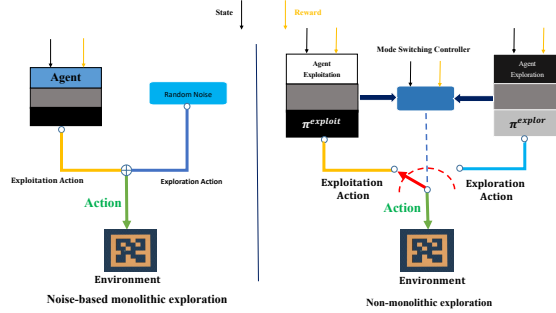


Fig. 1: The diagram illustrates noise-based monolithic exploration (left) and non-monolithic exploration (right). In noise-based monolithic exploration, an agent’s action and noise serve roles in exploitation and exploration, respectively. ‘Action’ represents the final action taken by the agent in the environment, resulting from the addition of the action and noise. In non-monolithic exploration, separate agents for exploitation and exploration operate independently, facilitated by a mode-switching controller. This controller assesses the state of either one agent or both agents, allowing them to pursue their respective objectives.

- Blind switching : not considering any state
- Informed switching : the opposite of blind switching and switching informed by the agent’s internal state
- Value promise trigger : triggering according to Eq. 1
- Starting mode : an episode started in explore mode or in exploit mode

A.4 Flexibility to the exploration process

- Bandit adaptation : a meta-controller parameterised by termination probability or target rate
- Homeostasis : the binary switching mechanism based on Algorithm 1

A.5 Homeostasis

The Algorithm 1 of Homeostasis is used for the evaluation implementation of our model. A sequence of scalar signals $x_t \in \mathbb{R}$ is transformed to a sequence of binary switching decisions $y_t \in \{0, 1\}$ for $1 \leq t \leq T$. The binary switching decisions y_t are used for target rate ρ which is the average number of switches approximated to $\frac{1}{T} \sum_t y_t \approx \rho \in \{0.1, 0.01, 0.001, 0.0001\}$.

A.6 Variants of intra-episodic exploration

Here is a classification for two-mode intra-episodic exploration for non-monolithic exploration model.

- Explore mode : 1. uniform random (XU). 2. intrinsic reward (XI).

Table 1: The tasks of D4RL environment for experiments.

Task	D4RL Environment Name
antmaze-umaze	antmaze-umaze-v2
antmaze-umaze-diverse	antmaze-umaze-diverse-v2
antmaze-medium-play	antmaze-medium-play-v2
antmaze-medium-diverse	antmaze-medium-diverse-v2
antmaze-large-play	antmaze-large-play-v2
antmaze-large-diverse	antmaze-large-diverse-v2
halfcheetah-random	halfcheetah-random-v2
hopper-random	hopper-random-v2
walker-random	walker-random-v2
halfcheetah-medium	halfcheetah-medium-v2
hopper-medium	hopper-medium-v2
walker-medium	walker-medium-v2
halfcheetah-medium-replay	halfcheetah-medium-replay-v2
hopper-medium-replay	hopper-medium-replay-v2
walker-medium-replay	walker-medium-replay-v2

- Explore duration : 1. fixed number of steps (1, 10, 100). 2. adaptively picked by a bandit (represented by *). 3. symmetric switching between entering and exiting explore mode (represented by =).
- Trigger type : 1. blind. 2. informed (based on value promise).
- Exploit duration : 1. blind triggers : a. fixed number of steps (10, 100, 1000, 10000), indirectly defined by a probability of terminating (0.1, 0.01, 0.001, 0.0001) (represented by n*). b. adaptively picked by a bandit over these choices (represented by p*). 2. informed triggers : a. indirectly parameterised by a target rate in (0.1, 0.01, 0.001, 0.0001). b. a bandit over them (represented by p*). 'Informed triggers' (a or b) is transformed into an adaptive switching threshold by homeostasis.
- Starting mode : 1. greedy (G). 2. explore (X).

XU-intra(100,informed,p*,X) is an example instance based on the above classification for the implementation.

Algorithm 1 Homeostasis algorithm taken from [1]

```

1: Require:
   Target rate  $\rho$ 
2: Initialize:
    $\bar{x} \leftarrow 0, \bar{x}^2 \leftarrow 1, x^+ \leftarrow 1$ 
3: for  $t \in \{1, \dots, T\}$  do
4:   obtain next scalar signal return  $x_t$ 
5:   set time-scale  $\tau \leftarrow \min(t, \frac{100}{\rho})$ 
6:   update moving average  $\bar{x} \leftarrow (1 - \frac{1}{\tau})\bar{x} + \frac{1}{\tau}x_t$ 
7:   update moving variance  $\bar{x}^2 \leftarrow (1 - \frac{1}{\tau})\bar{x}^2 + \frac{1}{\tau}(x_t - \bar{x})^2$ 
8:   standardise and exponentiate  $x^+ \leftarrow \exp\left(\frac{x_t - \bar{x}}{\sqrt{\bar{x}^2}}\right)$ 
9:   update transformed moving average
10:   $\bar{x}^+ \leftarrow (1 - \frac{1}{\tau})\bar{x}^+ + \frac{1}{\tau}x^+$ 
11:  sample  $y_t \sim \text{Bernoulli}\left(\min\left(1, \rho \frac{\bar{x}^+}{x^+}\right)\right)$ 
12: end for

```

Table 2: Hyper-parameters used in the implementation.

Hyper-parameters	Value
number of parallel env	1
discount	0.99
replay buffer size	1e6
batch size	256
MLP hidden layer size	[256, 256]
learning rate	3e-4
initial collection steps	5000
target update speed	5e-3
expectile value τ	0.9 (0.7)
inverse temperature α^{-1}	10(3)
number of offline iterations	1M
number of online iterations	1M
number of iteration per rollout step	1
target entropy (SAC)	-d

B Tasks and Hyper-Parameters

Task and hyper-parameters are adopted from PEX without modification. The following Table 1 and 2 represents the tasks of D4RL environment and the hyper-parameters for experiments, respectively.

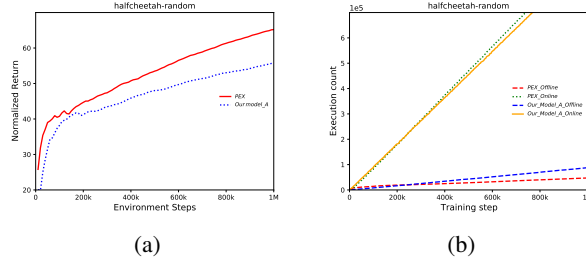


Fig. 2: Normalized return (left) and execution count (right) of PEX and our model (to stress an arbitrary agent referred to as *Ourmodel_A*) on ‘halfcheetah-random’ beyond the execution count of PEX. In (b), the execution counts of the offline policy and online policy of PEX or our model are referred to as *PEX_Offline* and *PEX_Online* or *Our_Model_A_Offline* and *Our_Model_A_Online*, respectively.

C Ablation study

According to the experiments of our model, its performance is certainly inferior to that of PEX, as shown in Fig. 2, particularly when the execution count of each policy in our model exceeds that of PEX. The following illustrates the condition using the notation employed in the manuscript.

$$\begin{aligned}
 &OurModel_Online \leq PEX_Online \\
 &\quad \text{and} \\
 &OurModel_Offline \geq PEX_Offline
 \end{aligned} \tag{2}$$

This ‘Ablation study’ indicates that our model relies less on offline policy and more on online policy than PEX. Finally, it shows that the execution count of each policy relying on the adjustment of three key parameters— ρ , ‘*explore_fixed_steps*’, and ‘*update_timestep*’—has an impact on the performance of our model.

References

1. Pislari, M., Szepesvari, D., Ostrovski, G., Borsa, D., Schaul, T.: When should agents explore? arXiv preprint arXiv:2108.11811 (2021)

Table 3: Key Notations.

Symbol	Meaning
t	action step
s	current state
s'	next state
r	reward
a	The final action to an environment
π_β or π^{off}	The policy trained from offline training stage
π_θ or π^{on}	The policy trained from online training stage
Q^{off}	The state-action value function obtained from offline training stage
Π	The policy set composing of $[\pi_\beta, \pi_\theta]$ or $[\pi^{\text{off}}, \pi^{\text{on}}]$
$\tilde{\pi}$	The single composite policy formed from policy set Π
<i>Action</i>	The action of $\tilde{\pi}$ to an environment
a_t^{off}	The action of π^{off}
a_t^{on}	The action of π^{on}
f_{homeo}	The function of homeostasis of the policy of Exploit, π^{off}
ρ	The preset value of target rate, i.e. the average number of switches of the reference model, which is selected among $[0.0001, 0.9]$ for fine-training. See more in ‘A.5 Homeostasis’ of ‘supplementary document.pdf’.
D_{promise}	The value promise discrepancy of π^{off}
D_{off}	Offline replay buffer
D_{on}	Online replay buffer
$D_{\text{Down,OPT}}$	The optimal replay buffer of a downstream task