

Supplementary Information

April 25, 2024

1 algorithms

Algorithm 1 Fix-weight

```
1: Input: State  $S_t$ , weights  $\mathbf{w}$ 
2: Initialize:  $K$  weight-conditioned Actor-critic networks  $\pi_i$  and  $v^{\pi_i}$ , predetermined sub-space  $\mathbf{W}_i$ , and memory buffer  $\mathbf{E}$  size of  $B$ .
3: for  $i = 0$  to  $n$  do
4:   for  $t = 1$  to  $B$  do
5:      $\mathbf{w} \leftarrow \text{get pivot weight}(\mathbf{W}_i)$ 
6:      $a_t \leftarrow \pi_i(S_t, \mathbf{w})$ 
7:      $S_{t+1}, \mathbf{r}_t \leftarrow \text{simulator}(a_t)$ 
8:      $\mathbf{E} \leftarrow \mathbf{E} \cup \langle S_t, a_t, \mathbf{w}, \mathbf{r}_t, S_{t+1} \rangle$ 
9:      $S_t \leftarrow S_{t+1}$ 
10:   end for
11:   sample  $\langle S_t, a_t, \mathbf{w}, \mathbf{r}_t, S_{t+1} \rangle \leftarrow \mathbf{E}$ 
12:    $\theta \leftarrow \theta + \eta \nabla_{\theta} \log \pi_{\theta}(s, a; \mathbf{w}) A^{\pi}(s_t, a_t; \mathbf{w})$ 
13:    $\phi \leftarrow \phi + ||\mathbf{V}^{\pi_i}(S; \mathbf{w}) - \mathbf{V}^{\pi_i}(S_{t+1}; \mathbf{w})||^2$ 
14:   clear  $\mathbf{E}$ 
15: end for
```

Algorithm 2 Random-weight

```
1: Input: State  $S_t$ , weights  $\mathbf{w}$ 
2: Initialize:  $K$  weight-conditioned Actor-critic networks  $\pi_i$  and  $v^{\pi_i}$ , predetermined sub-space  $\mathbf{W}_i$ , weight change period  $k$ , and memory buffer  $\mathbf{E}$  size of  $B$ .
3: for  $i = 0$  to  $K$  do
4:   for  $t = 0$  to  $B$  do
5:     if  $t \bmod k = 0$  then
6:        $\mathbf{w}_t \leftarrow \text{uniform sample}(\mathbf{W}_i)$ 
7:     end if
8:      $a_t \leftarrow \pi_i(S_t, \mathbf{w}_t)$ 
9:      $S_{t+1}, \mathbf{r}_t \leftarrow \text{simulator}(a_t)$ 
10:     $\mathbf{E} \leftarrow \mathbf{E} \cup \langle S_t, a_t, \mathbf{w}_t, \mathbf{r}_t, S_{t+1} \rangle$ 
11:     $S_t \leftarrow S_{t+1}$ 
12:  end for
13:  sample  $\langle S_t, a_t, \mathbf{w}_t, \mathbf{r}_t, S_{t+1} \rangle \leftarrow \mathbf{E}$ 
14:   $\theta \leftarrow \theta + \eta \nabla_{\theta} \log \pi_{\theta}(s, a; \mathbf{w}_t) A^{\pi}(s_t, a_t; \mathbf{w}_t)$ 
15:   $\phi \leftarrow \phi + ||\mathbf{V}^{\pi_i}(S; \mathbf{w}_t) - \mathbf{V}^{\pi_i}(S_{t+1}; \mathbf{w}_t)||^2$ 
16:  clear  $\mathbf{E}$ 
17: end for
```

Algorithm 3 UCB-MOPPO

```
1: Input: State  $S_t$ , weights  $\mathbf{w}$ 
2: Initialize:  $K$  weight conditioned Actor-critic network  $\pi_i$  and  $v^{\pi_i}$ , predetermined sub-space  $\mathbf{W}_i$  size of  $M$ , current working weight space  $\widetilde{\mathbf{W}}_i$  size of  $M$ , each objective has  $m$  dimensions, warm-up iterations  $Q$ , objective value collection interval in every  $C$  iterations, total  $K \times m$  surrogate model  $\psi_\theta$ , and dynamic weight experience pool  $\mathcal{E}$  size of  $B$ .
3: for  $t = 0$  to  $T$  do
4:   ► Warm-up Stage
5:   for  $i = 0$  to  $n$  do
6:      $\widetilde{\mathbf{W}}_i \leftarrow$  get pivot weights( $\mathbf{W}_i$ )
7:      $\pi_i^* \leftarrow$  Fix Weight Optimisation( $\pi_i, \widetilde{\mathbf{W}}_i$ )
8:   end for
9:   ► Collect objective value from simulator
10:  if  $t \bmod C = 0$  then
11:    for  $i = 0$  to  $n$  do
12:      for  $\mathbf{w}$  in  $\widetilde{\mathbf{W}}_i$  do
13:         $V^{\pi_i} \leftarrow$  simulator( $\pi_i, \mathbf{w}$ )
14:      end for
15:    end for
16:  end if
17:  if  $t > Q$  then
18:    ► Construct Training Data for Prediction Model:
19:    for  $i = 0$  to  $n$  do
20:      for  $\mathbf{w}$  in  $\widetilde{\mathbf{W}}_i$  do
21:        for  $k = 1$  to  $\frac{Q}{C}$  do
22:           $\delta V_k^{\pi_i} \leftarrow V_k^{\pi_i} - V_{k-1}^{\pi_i}$ 
23:        end for
24:      end for
25:    end for
26:    ► Update Prediction Model:
27:    for  $i = 0$  to  $n$  do
28:      for  $j = 0$  to  $m$  do
29:        for  $\mathbf{w}$  in  $\widetilde{\mathbf{W}}_i$  do
30:           $\theta \leftarrow \theta +$  grid search( $\psi_{\theta_j}^{\pi_i}, \delta V_k^{\pi_i}$ )
31:        end for
32:      end for
33:      ► Preference Search:
34:      for  $\mathbf{w}$  in  $\mathbf{W}_i$  do
35:        for  $\mathbf{w}$  in  $\mathbf{W}_i$  do
36:           $V_{\mathbf{w}}^{\pi_i} \leftarrow$  simulator( $\pi_i, \mathbf{w}$ )
37:           $L \leftarrow$  append( $V_{\mathbf{w}}^{\pi_i}$ )
38:        end for
39:        for  $j = 0$  to  $m$  do
40:           $\delta \hat{V}_{\mathbf{w},j}^{\pi_i} \leftarrow \psi_{\theta_j}^{\pi_i}(\mathbf{w}), \mathbf{w} \in \mathbf{W}_i$ 
41:           $\delta \tilde{V}_{\mathbf{w},j}^{\pi_i} \leftarrow \mu(\delta \hat{V}_{\mathbf{w},j}^{\pi_i}) + \beta(t) * \sigma(\delta \hat{V}_{\mathbf{w},j}^{\pi_i})$ 
42:           $\tilde{V}_{\mathbf{w},j}^{\pi_i} \leftarrow V_{\mathbf{w},j}^{\pi_i} + \delta \tilde{V}_{\mathbf{w},j}^{\pi_i}$ 
43:        end for
44:         $F \leftarrow \{\tilde{V}_{\mathbf{w}}^{\pi_i}\} \cup L \setminus \{V_{\mathbf{w}}^{\pi_i}\}$ 
45:         $\mathcal{D} \leftarrow$  append( $\langle \mathbf{w}_i, \mathcal{H}(F) \rangle$ )
46:      end for
47:      ► Update Working Preference Pool:
48:       $\{\mathbf{w}_i, i \in (0, M)\} \leftarrow$  Sort  $\omega$  of  $\mathcal{D}$  by  $\mathcal{H}(F)$  in descending order
49:       $\widetilde{\mathbf{W}}_i \leftarrow \{\mathbf{w}_i, i \in (0, M)\}$ 
50:    end for
51:  end if
52: end for
```

2 Benchmark Problems

This section provide detail objective return for each problems. Where C in the following equations is live bonus.

2.1 Swimmer-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^8, \mathcal{A} \in \mathbb{R}^2$.

The first objective is forward speed in x axis:

$$R_1 = v_x \quad (1)$$

The second objective is energy efficiency:

$$R_2 = 0.3 - 0.15 \sum_i a_i^2, \quad a_i \in (-1, 1) \quad (2)$$

2.2 HalfCheetah-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{17}, \mathcal{A} \in \mathbb{R}^6$.

The first objective is forward speed in x axis:

$$R_1 = \min(v_x, 4) + C \quad (3)$$

The second objective is energy efficiency:

$$R_2 = 4 - \sum_i a_i^2 + C, \quad a_i \in (-1, 1) \quad (4)$$

$$C = 1 \quad (5)$$

2.3 Walker2d-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{17}, \mathcal{A} \in \mathbb{R}^6$.

The first objective is forward speed in x axis:

$$R_1 = v_x + C \quad (6)$$

The second objective is energy efficiency:

$$R_2 = 4 - \sum_i a_i^2 + C, \quad a_i \in (-1, 1) \quad (7)$$

$$C = 1 \quad (8)$$

2.4 Ant-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{27}, \mathcal{A} \in \mathbb{R}^8$.

The first objective is forward speed in x axis:

$$R_1 = v_x + C \quad (9)$$

The second objective is forward in y axis:

$$R_2 = v_y + C \quad (10)$$

$$C = 1 - 0.5 \sum_i a_i^2, \quad a_i \in (-1, 1) \quad (11)$$

2.5 Hopper-v2

Observation and action space: $\mathcal{S} \in \mathbb{R}^{11}, \mathcal{A} \in \mathbb{R}^3$.

The first objective is forward speed in x axis:

$$R_1 = 1.5v_x + C \quad (12)$$

The second objective is jumping height:

$$R_2 = 12(h - h_{init}) + C \quad (13)$$

$$C = 1 - 2e^{-4} \sum_i a_i^2, \quad a_i \in (-1, 1) \quad (14)$$

2.6 Hopper-v3

Observation and action space: $\mathcal{S} \in \mathbb{R}^{11}, \mathcal{A} \in \mathbb{R}^3$.

The first objective is forward speed in x axis:

$$R_1 = 1.5v_x + C \quad (15)$$

The second objective is jumping height:

$$R_2 = 12(h - h_{init}) + C \quad (16)$$

The third objective is energy efficiency:

$$R_3 = 4 - \sum_i a_i^2 + C \quad (17)$$

$$C = 1 \quad (18)$$