

Analysis of Markov Decision Processes under Parameter Uncertainty Online Companion

Peter Buchholz, Iryna Dohndorf, and Dimitri Scheftelowitsch

Department of Computer Science, TU Dortmund
{peter.buchholz,iryana.dohndorf,dimitri.scheftelowitsch}@cs.tu-dortmund.de

Abstract. Markov Decision Processes (MDPs) are a popular decision model for stochastic systems. Introducing uncertainty in the transition probability distribution by giving upper and lower bounds for the transition probabilities yields the model of Bounded Parameter MDPs (BMDPs) which captures many practical situations with limited knowledge about a system or its environment. In this paper the class of BMDPs is extended to Bounded Parameter Semi Markov Decision Processes (BSMDPs). The main focus of the paper is on the introduction and numerical comparison of different algorithms to compute optimal policies for BMDPs and BSMDPs; specifically, we introduce and compare variants of value and policy iteration.

The paper delivers an empirical comparison between different numerical algorithms for BMDPs and BSMDPs, with an emphasis on the required solution time.

Keywords: (Bounded Parameter) (Semi-)Markov Decision Process, Discounted Reward, Average Reward, Value Iteration, Policy Iteration

1 Algorithms

1.1 Solution methods for MDPs

Algorithm 1 Value iteration for discrete-time MDPs with discounted reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$, discount factor γ ;

- 1: Specify $\mathbf{v}^{(0)} \geq \mathbf{0}$, $\epsilon > 0$ and set $k = 0$;
 - 2: **repeat**
 - 3: **for** $i \in \mathcal{S}$ **do**
 - 4: $\mathbf{v}^{(k+1)}(i) = \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \gamma \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{v}^{(k)}(j) \right)$;
 - 5: $k = k + 1$;
 - 6: **until** $\left\| \mathbf{v}^{(k-1)} - \mathbf{v}^{(k)} \right\| < \epsilon \frac{1-\gamma}{2\gamma}$
 - 7: Choose $\pi(i) \in \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \gamma \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{v}^{(k)}(j) \right)$ for all $i \in \mathcal{S}$;
 - 8: **return** An ϵ -optimal policy π , value vector $\mathbf{v}^{(k)}$;
-

Algorithm 2 Value iteration for discrete-time MDPs with expected total reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$;

- 1: Specify $\mathbf{v}^{(0)} \geq \mathbf{0}$, $\epsilon > 0$ and set $k = 0$;
 - 2: **repeat**
 - 3: **for** $i \in \mathcal{S}$ **do**
 - 4: $\mathbf{v}^{(k+1)}(i) = \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{v}^{(k)}(j) \right)$;
 - 5: $k = k + 1$;
 - 6: **until** $\left\| \mathbf{v}^{(k-1)} - \mathbf{v}^{(k)} \right\| < \epsilon$
 - 7: $\pi(i) \in \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{v}^{(k)}(j) \right)$ for all $i \in \mathcal{S}$;
 - 8: **return** An ϵ -optimal policy π ;
-

Algorithm 3 Policy iteration for discrete-time MDPs with discounted reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$, discount factor γ

- 1: Specify $\pi^{(0)} \in \Pi$ some pure initial policy and set $k = 0$;
- 2: **repeat**
- 3: **(Policy evaluation)** Solve

$$\mathbf{r}^{\pi^{(k)}} = \left(\mathbf{I} - \gamma \mathbf{P}^{\pi^{(k)}} \right) \mathbf{v}^{(k)};$$

- 4: **(Policy improvement)** Choose $\pi^{(k+1)}$ to satisfy

$$\pi^{(k+1)} = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}^a + \gamma \mathbf{P}^a \mathbf{v}^{(k)} \right)$$

- 5: choosing $\pi^{(k+1)}(i) = \pi^{(k)}(i)$ when possible;
 - 5: $k = k + 1$;
 - 6: **until** $\pi^{(k)} = \pi^{(k-1)}$
 - 7: **return** An optimal policy $\pi^* = \pi^{(k)}$;
-

Algorithm 4 Policy iteration for discrete-time MDPs with expected total reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$;

1: Specify $\boldsymbol{\pi}^{(0)} \in \Pi$ some pure initial policy and set $k = 0$;

2: **repeat**

3: **(Policy evaluation)** Solve

$$\mathbf{r}^{\boldsymbol{\pi}^{(k)}} = (\mathbf{I} - \mathbf{P}^{\boldsymbol{\pi}^{(k)}}) \mathbf{v}^{(k)};$$

4: **(Policy improvement)** Choose $\boldsymbol{\pi}^{(k+1)}$ to satisfy

$$\boldsymbol{\pi}^{(k+1)} = \arg \max_{a \in \mathcal{A}} (\mathbf{r}^a + \mathbf{P}^a \mathbf{v}^{(k)})$$

choosing $\boldsymbol{\pi}^{(k+1)}(i) = \boldsymbol{\pi}^{(k)}(i)$ when possible;

5: $k = k + 1$;

6: **until** $\boldsymbol{\pi}^{(k-1)} = \boldsymbol{\pi}^{(k)}$

7: **return** An optimal policy $\boldsymbol{\pi}^* = \boldsymbol{\pi}^{(k)}$;

1.2 Solution methods for SMDPs

Algorithm 5 Relative value iteration for discrete-time MDPs with average reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$

- 1: Specify $\mathbf{v}^{(0)} \geq \mathbf{0}$, $\epsilon > 0$, set $k = 0$, and choose one state $i_0 \in \mathcal{S}$;
 - 2: **repeat**
 - 3: $\mathbf{w}^{(k)} = \mathbf{v}^{(k)} - \epsilon \mathbf{v}^{(k)}(i_0)$;
 - 4: **for** $i \in \mathcal{S}$ **do**
 - 5: $\mathbf{v}^{(k+1)}(i) = \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{w}^{(k)}(j) \right)$;
 - 6: $k = k + 1$;
 - 7: **until** $\max_{i \in \mathcal{S}} \left(\mathbf{v}^{(k+1)}(i) - \mathbf{v}^{(k)}(i) \right) - \min_{i \in \mathcal{S}} \left(\mathbf{v}^{(k+1)}(i) - \mathbf{v}^{(k)}(i) \right) < \epsilon$
 - 8: $\pi(i) \in \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}^a(i) + \sum_{j \in \mathcal{S}} \mathbf{P}^a(i, j) \mathbf{w}^{(k)}(j) \right)$ for all $i \in \mathcal{S}$;
 - 9: Set $\bar{\pi} = \pi$, $G = \mathbf{w}^{(k)}(1)$ and $\mathbf{h} = \mathbf{w}$;
 - 10: **return** An ϵ -optimal policy $\bar{\pi}$, average gain G and bias vector \mathbf{h} ;
-

Algorithm 6 Policy iteration for discrete-time MDPs with average reward criterion

Require: MDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$

- 1: Specify $\pi^{(0)} \in \Pi$ some pure initial policy and set $k = 0$;
- 2: **repeat**
- 3: **(Policy evaluation)** Solve

$$\mathbf{r}^{\pi^{(k)}} = \left(\mathbf{I} - \mathbf{P}^{\pi^{(k)}} \right) \bar{\mathbf{g}}^{(k)} + G \mathbf{1}$$

by setting $\bar{\mathbf{g}}(i_0) = 0$ for some fixed state $i_0 \in \mathcal{S}$. Compute $\mathbf{H}^{(k)} = \left(\mathbf{I} - \mathbf{P}^{\pi^{(k)}} \right)$.

Then $\bar{\mathbf{H}}^{(k)}$ is the matrix with the column corresponding to state i_0 replaced by a column of 1's. Solve the linear system

$$\mathbf{r}^{\pi^{(k)}} = \bar{\mathbf{H}}^{(k)} \mathbf{w}$$

where $G^{(k)}$ is the i_0 th component of the solution vector \mathbf{w} and $\mathbf{h}^{(k)}(i) = \mathbf{w}(i)$ for $i \neq i_0$;

- 4: **(Policy improvement)** Choose $\pi^{(k+1)}$ to satisfy

$$\pi^{(k+1)} = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}^a + \mathbf{P}^a \mathbf{h}^{(k)} \right)$$

choosing $\pi^{(k+1)}(i) = \pi^{(k)}(i)$ when possible;

- 5: $k = k + 1$;
 - 6: **until** $\pi^{(k-1)} = \pi^{(k)}$
 - 7: Set $\bar{\pi}^* = \pi^{(k)}$, $G^* = G^{(k-1)}$, $\mathbf{h} = \mathbf{h}^{(k-1)}$
 - 8: **return** An optimal policy $\bar{\pi}^*$, the optimal average gain G^* and the deviation vector \mathbf{h} ;
-

Algorithm 7 Uniformization method for SMDPs with average reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$, time vectors $\{\mathbf{y}^a\}_{a \in \mathcal{A}}$ (average sojourn times in states)

- 1: Choose $\eta = \min_{i \in \mathcal{S}} \min_{a \in \mathcal{A}} \mathbf{y}^a(i) / (1 - \mathbf{P}^a(i, i))$;
 - 2: **for** $a \in \mathcal{A}$ **do**
 - 3: $\bar{\mathbf{s}}^a(i) = \mathbf{r}^a(i) / \mathbf{y}^a(i)$;
 - 4: **for** $a \in \mathcal{A}$ **do**
 - 5: **for** $i \in \mathcal{S}$ **do**
 - 6: **for** $j \in \mathcal{S}$ **do**
 - 7: **if** $i \neq j$ **then**
 - 8: $\bar{Q}^a(i, j) = \eta \frac{\mathbf{P}^a(i, j)}{\mathbf{y}^a(i)}$;
 - 9: **else**
 - 10: $\bar{Q}^a(i, j) = 1 + \eta \frac{\mathbf{P}^a(i, j) - 1}{\mathbf{y}^a(i)}$;
 - 11: **return** Discrete-time MDP $(\mathcal{S}, \mathcal{A}, (\bar{Q}^a)_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}})$, η ;
-

Algorithm 8 Value iteration for discrete-time SMDPs with average reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$, time vectors $\{\mathbf{y}^a\}_{a \in \mathcal{A}}$ (average sojourn times in states)

- 1: Apply Algorithm 7 to transform the SMDP in an according to the average reward equivalent MDP $\{\mathcal{S}, \mathcal{A}, \{\bar{Q}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}}\}$. Save η ;
 - 2: Use value iteration Algorithm 5 to analyze the MDP;
 - 3: **return** An ϵ -optimal policy $\bar{\pi}$, average gain G and bias vector $\mathbf{h} = \eta \mathbf{h}$;
-

Algorithm 9 Policy iteration for discrete-time SMDPs with average reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, \mathbf{p})$, time vectors $\{\mathbf{y}^a\}_{a \in \mathcal{A}}$ (average sojourn times in states)

- 1: Apply Algorithm 7 to transform the SMDP in an according to the average reward equivalent MDP $\{\mathcal{S}, \mathcal{A}, \{\bar{Q}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}}\}$. Save η ;
 - 2: Use policy iteration Algorithm 6 to analyze the MDP;
 - 3: **return** An optimal policy $\bar{\pi}^*$, average gain G^* and bias vector $\mathbf{h} = \eta \mathbf{h}$;
-

Algorithm 10 Transformation method for SMDPs with discounted reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}^a)_{a \in \mathcal{A}}, (\mathbf{r}^a)_{a \in \mathcal{A}}, ((\mathbf{p}^{(a,i)}, \mathbf{D}_0^{(a,i)})_{a \in \mathcal{A}, i \in \mathcal{S}})$, discount factor β ;

- 1: **for** $a \in \mathcal{A}$ **do**
- 2: Let $\{(\mathbf{p}^{(i)}, \mathbf{D}_0^{(i)})\}_{i \in \mathcal{S}}$ be the set of Phase-type distributions corresponding to the action a ;
- 3: **for** $i \in \mathcal{S}$ **do**
- 4: Compute $\mathbf{s}^a(i) = \mathbf{r}^a(i) \int_0^\infty (1 - F^a(i, t)) e^{-\beta t} dt$ and for all $j \in \mathcal{S}$

$$\bar{\mathbf{Q}}^a(i, j) = \mathbf{P}^a(i, j) \int_0^\infty f^a(i, t) e^{-\beta t} dt$$
with the uniformization based method [1] using the following data.

Compute $\mathbf{d}_1^{(i)} = -\mathbf{D}_0^{(i)} \mathbf{I}$;
Set $\mathbf{P}^{(i)} = \mathbf{D}_0^{(i)} - \beta \mathbf{I}$ and $\lambda = \max_{i, j \in \mathcal{S}} |\mathbf{P}^{(i)}(i, j)|$;
Compute $\mathbf{P}^{(i)} = \frac{1}{\lambda} \mathbf{P}^{(i)} + \mathbf{I}$, $\mathbf{d}_1^{(i)} = \frac{1}{\lambda} \mathbf{d}_1^{(i)}$ and the time step $\Delta = 1/\lambda$;
- 5: **return** Discrete-time discounted MDP $(\mathcal{S}, \mathcal{A}, (\bar{\mathbf{Q}}^a)_{a \in \mathcal{A}}, (\mathbf{s}^a)_{a \in \mathcal{A}})$;

Algorithm 11 Value iteration for discrete-time SMDPs with discounted reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, \{\mathbf{P}^a\}_{a \in \mathcal{A}}, \{\mathbf{r}^a\}_{a \in \mathcal{A}}, \{(\mathbf{p}^{(a,i)}, \mathbf{D}_0^{(a,i)})\}_{a \in \mathcal{A}, i \in \mathcal{S}})$, discount factor β ;

- 1: Apply transformation Algorithm 10 to transform the SMDP in an according to the discounted reward equivalent MDP $(\mathcal{S}, \mathcal{A}, \{\bar{\mathbf{Q}}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}})$;
- 2: Use value iteration Algorithm 2 to analyze the MDP $\{\mathcal{S}, \mathcal{A}, \{\bar{\mathbf{Q}}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}}\}$ according to the expected total reward criterion;
- 3: **return** An ϵ -optimal policy $\boldsymbol{\pi}$;

Algorithm 12 Policy iteration for discrete-time SMDPs with discounted reward criterion

Require: SMDP $(\mathcal{S}, \mathcal{A}, \{\mathbf{P}^a\}_{a \in \mathcal{A}}, \{\mathbf{r}^a\}_{a \in \mathcal{A}}, \{(\mathbf{p}^{(a,i)}, \mathbf{D}_0^{(a,i)})\}_{a \in \mathcal{A}, i \in \mathcal{S}})$, discount factor β ;

- 1: Apply transformation Algorithm 10 to transform the SMDP in an according to the discounted reward equivalent MDP $(\mathcal{S}, \mathcal{A}, \{\bar{\mathbf{Q}}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}})$;
- 2: Use policy iteration Algorithm 4 to analyze the MDP $\{\mathcal{S}, \mathcal{A}, \{\bar{\mathbf{Q}}^a\}_{a \in \mathcal{A}}, \{\mathbf{s}^a\}_{a \in \mathcal{A}}\}$ according to the expected total reward criterion;
- 3: **return** Optimal policy $\boldsymbol{\pi}^*$;

1.3 Solution methods for BMDPs

Algorithm 13 Interval value iteration for discrete-time BMDPs with discounted reward criterion

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$, discount factor $(\gamma_{\downarrow}^a)_{a \in \mathcal{A}}$, *pessimistic* is *true* when the optimal lower bound has to be computed and *false* when the optimal upper bound has to be computed;

- 1: Specify $\mathbf{v}^{(0)} \geq \mathbf{0}$, $\boldsymbol{\pi}^{(0)} \geq \mathbf{0}$, $\epsilon > 0$ and set $k = 0$;
- 2: $\rightarrow = \downarrow$ if *pessimistic*, otherwise $\rightarrow = \uparrow$
- 3: $\gamma^* = \max \{ \gamma_{\rightarrow}^a(i) \mid (i, a) \in \mathcal{S} \times \mathcal{A} \}$
- 4: **repeat**
- 5: **for** $i \in \mathcal{S}$ **do**
- 6: $[\mathbf{v}^{(k+1)}(i), \boldsymbol{\pi}^{(k+1)}(i)] = \text{interval_value}(i, \boldsymbol{\pi}^{(k)}(i), (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}}, (\gamma_{\rightarrow}^a)_{a \in \mathcal{A}}, \mathbf{v}^{(k)}, \epsilon, \textit{pessimistic})$;
- 7: $k = k + 1$;
- 8: **until** $\| \mathbf{v}^{(k+1)} - \mathbf{v}^{(k)} \| < \epsilon \frac{1 - \gamma^*}{2\gamma^*}$
- 9: **return** An ϵ -optimal policy $\boldsymbol{\pi}^{(k)}$, value vector $\mathbf{v}^{(k)}$;

```

1: function INTERVAL_VALUE(state  $i$ , current decision  $ai$ ,  $(P_{\uparrow}^a)_{a \in \mathcal{A}}$ ,  $(r_{\uparrow}^a)_{a \in \mathcal{A}}$ ,  $(\gamma^a)_{a \in \mathcal{A}}$ ,
    $\mathbf{v}$ ,  $\epsilon$ , pessimistic)
2:    $w = -1.0e + 12$ ;
3:   if pessimistic then
4:      $(i_1, i_2, \dots, i_n) \leftarrow$  ascending order of states with respect to states' values  $\mathbf{v}$ ;
5:   else
6:      $(i_1, i_2, \dots, i_n) \leftarrow$  descending order of states with respect to states' values  $\mathbf{v}$ ;
7:    $\mathbf{p} = P_{\downarrow}^a$ ;
8:   for  $a \in \mathcal{A}$  do
9:      $val = \gamma \mathbf{p} \mathbf{v}$  ;
10:    if pessimistic then
11:       $val = val + r_{\downarrow}^a(i)$  ;
12:    else
13:       $val = val + r_{\uparrow}^a(i)$  ;
14:     $r = \mathbf{p}\mathbf{1}$ ;
15:    for  $j \in (i_1, i_2, \dots, i_n)$  do
16:      if  $P_{\uparrow}^a(i, j) > P_{\downarrow}^a(i, j)$  ; then
17:         $m = \min(P_{\uparrow}^a(i, j) - \mathbf{p}(j), 1 - r)$ ;
18:         $val = val + \gamma^a(i) \cdot m\mathbf{v}(j)$ ;
19:         $r = r + m$ ;
20:      if  $r \geq 1 - 10\epsilon$ ; then return ;
21:    if  $val > w$  then
22:       $w = val$ ;
23:       $ai = a$ ;
24:  return  $ai, w$ ;

```

Algorithm 14 Policy iteration 1 for discrete-time BMDPs with discounted reward criterion

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$, discount factor $(\gamma_{\downarrow}^a)_{a \in \mathcal{A}}$, *pessimistic* is *true* when the optimal lower bound has to be computed and *false* when the optimal upper bound has to be computed;

1: Specify $\phi^{(1)} \in \Pi$ some pure initial policy, $\mathbf{v}^{(0)} = \mathbf{r}^{\phi^{(1)}}$ and set $k = 1$;

2: **if** *pessimistic* **then**

3: $\Gamma = \text{diag}(\gamma_{\downarrow}^{\phi^{(k)}})$;

$$\mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(k-1)}) = \arg \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\phi^{(k)}}} (\Gamma \mathbf{P} \mathbf{v}^{(k-1)});$$

Solve

$$\mathbf{r}_{\downarrow}^{\phi^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(k-1)}) \right) \mathbf{v}^{(k)};$$

4: **for** $i \in \mathcal{S}$ **do**

5: **for** $a \in \mathcal{A}$ **do**

6: $f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}^{(k)}) = \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\Gamma \mathbf{P}(i \bullet) \mathbf{v}^{(k)})$;

7: Choose $\phi^{(k+1)}(i)$ to satisfy

$$\phi^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\downarrow}^a(i) + \gamma_{\downarrow}^a(i) f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}^{(k)}) \right)$$

keeping $\phi^{(k+1)}(i) = \phi^{(k)}(i)$ when possible;

8: **if** $\phi^{(k+1)} = \phi^{(k)}$ **then**

9: Set $\phi_{\downarrow}^* = \phi^{(k+1)}$ and terminate. Otherwise set $k = k + 1$ and go to Step 3;

10: **else**

11: $\Gamma = \text{diag}(\gamma_{\uparrow}^{\phi^{(k)}})$;

$$\mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(k-1)}) = \arg \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\phi^{(k)}}} (\Gamma \mathbf{P} \mathbf{v}^{(k-1)});$$

Solve

$$\mathbf{r}_{\uparrow}^{\phi^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(k-1)}) \right) \mathbf{v}^{(k)};$$

12: **for** $i \in \mathcal{S}$ **do**

13: **for** $a \in \mathcal{A}$ **do**

14: $f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}^{(k)}) = \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \mathbf{v}^{(k)})$;

15: Choose $\phi^{(k+1)}(i)$ to satisfy

$$\phi^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\uparrow}^a(i) + \gamma_{\uparrow}^a(i) f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}^{(k)}) \right)$$

keeping $\phi^{(k+1)}(i) = \phi^{(k)}(i)$ when possible;

16: **if** $\phi^{(k+1)} = \phi^{(k)}$ **then**

17: Set $\phi_{\uparrow}^* = \phi^{(k+1)}$ and terminate. Otherwise set $k = k + 1$ and go to Step 11;

18: **return** An optimal policy ϕ_{\downarrow}^* if *pessimistic* is *true* and ϕ_{\uparrow}^* if *pessimistic* is *false*;

Algorithm 15 Policy iteration 2 for discrete-time BMDPs with discounted reward criterion

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$, discount factor $(\gamma_{\downarrow}^a)_{a \in \mathcal{A}}$, *pessimistic* is *true* when the optimal lower bound has to be computed and *false* when the optimal upper bound has to be computed;

1: Specify $\phi^{(1)} \in \Pi$ some pure initial policy, $\mathbf{v}^{(0)} = \mathbf{r}^{\phi^{(1)}}$, $\epsilon > 0$ and set $k = 1$, $l = 1$;

2: **if** *pessimistic* **then**

3: **repeat**

4: $\Gamma = \text{diag}(\gamma_{\downarrow}^{\phi^{(k)}})$;

$M_{\downarrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(l-1)}) = \arg \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\phi^{(k)}}} (\Gamma \mathbf{P} \mathbf{v}^{(l-1)})$;

 Solve

$$\mathbf{r}_{\downarrow}^{\phi^{(k)}} = \left(\mathbf{I} - M_{\downarrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(l-1)}) \right) \mathbf{v}^{(l)};$$

5: **until** $\|\mathbf{v}^{(l)} - \mathbf{v}^{(l-1)}\| < \epsilon$;

6: $\mathbf{v} = \mathbf{v}^{(l)}$;

7: **for** $i \in \mathcal{S}$ **do**

8: **for** $a \in \mathcal{A}$ **do**

9: $f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}) = \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \mathbf{v})$;

10: Choose $\phi^{(k+1)}(i)$ to satisfy

$$\phi^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\downarrow}^a(i) + \gamma_{\downarrow}^a(i) f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}) \right)$$

 keeping $\phi^{(k+1)}(i) = \phi^{(k)}(i)$ when possible;

11: **if** $\phi^{(k+1)} = \phi^{(k)}$ **then**

12: Set $\phi_{\downarrow}^* = \phi^{(k+1)}$ and terminate. Otherwise set $k = k + 1$, $l = 1$, $\mathbf{v}^{(0)} = \mathbf{v}$ and go to Step 3;

13: **else**

14: **repeat**

15: $\Gamma = \text{diag}(\gamma_{\downarrow}^{\phi^{(k)}})$;

$M_{\uparrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(l-1)}) = \arg \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\phi^{(k)}}} (\Gamma \mathbf{P} \mathbf{v}^{(l-1)})$;

 Solve

$$\mathbf{r}_{\uparrow}^{\phi^{(k)}} = \left(\mathbf{I} - \gamma M_{\uparrow}(\mathbf{P}_{\downarrow}^{\phi^{(k)}}, \mathbf{v}^{(l-1)}) \right) \mathbf{v}^{(l)};$$

16: **until** $\|\mathbf{v}^{(l)} - \mathbf{v}^{(l-1)}\| < \epsilon$;

17: $\mathbf{v} = \mathbf{v}^{(l)}$;

18: **for** $i \in \mathcal{S}$ **do**

19: **for** $a \in \mathcal{A}$ **do**

20: $f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}) = \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \mathbf{v})$;

21: Choose $\phi^{(k+1)}(i)$ to satisfy

$$\phi^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\uparrow}^a(i) + \gamma_{\uparrow}^a(i) f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \mathbf{v}) \right)$$

 keeping $\phi^{(k+1)}(i) = \phi^{(k)}(i)$ when possible;

22: **if** $\phi^{(k+1)} = \phi^{(k)}$ **then**

23: Set $\phi_{\uparrow}^* = \phi^{(k+1)}$ and terminate. Otherwise set $k = k + 1$, $l = 1$, $\mathbf{v}^{(0)} = \mathbf{v}$ and go to Step 14;

24: **return** An optimal policy ϕ_{\downarrow}^* if *pessimistic* is *true* and ϕ_{\uparrow}^* if *pessimistic* is *false*;

Algorithm 16 Interval value iteration for discrete-time BMDPs with average reward criterion

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$, *pessimistic* is *true* when the optimal lower bound has to be computed and *false* when the optimal upper bound has to be computed;

- 1: Specify $\mathbf{v}^{(0)} \geq \mathbf{0}$, $\boldsymbol{\pi}^{(0)} \geq \mathbf{0}$, $\epsilon > 0$, set $k = 0$, and choose one state $i_0 \in \mathcal{S}$;
- 2: **repeat**
- 3: $\mathbf{w}^{(k)} = \mathbf{v}^{(k)} - \epsilon \mathbf{v}^{(k)}(i_0)$;
- 4: **for** $i \in \mathcal{S}$ **do**
- 5: $[\mathbf{v}^{(k+1)}(i), \boldsymbol{\pi}^{(k+1)}(i)] = \text{interval_value}(i, \boldsymbol{\pi}^{(k)}(i), (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}}, \mathbb{I}, \mathbf{w}^{(k)}, \epsilon, \text{pessimistic})$;
- 6: **until** $\max_{i \in \mathcal{S}} (\mathbf{v}^{(k+1)}(i) - \mathbf{v}^{(k)}(i)) - \min_{i \in \mathcal{S}} (\mathbf{v}^{(k+1)}(i) - \mathbf{v}^{(k)}(i)) < \epsilon$
- 7: Set $\bar{\boldsymbol{\pi}} = \boldsymbol{\pi}^{(k)}$, $G = \mathbf{w}^{(k)}(1)$ and $\mathbf{h} = \mathbf{w}$;
- 8: **return** An ϵ -optimal policy $\bar{\boldsymbol{\pi}}$, average gain G and bias vector \mathbf{h} ;

Algorithm 17 Policy iteration 1 for discrete-time BMDPs with average reward criterion when optimal *lower* bound should be computed

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$;

- 1: Specify $\bar{\boldsymbol{\phi}}^{(1)} \in \Pi$ some pure initial policy, $\bar{\mathbf{h}}^{(0)} = \mathbf{r}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(1)}}$ and set $k = 1$;
- 2: $\mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) = \arg \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}}} (\mathbf{P} \bar{\mathbf{h}}^{(k-1)})$;

Solve

$$\mathbf{r}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) \right) \bar{\mathbf{h}}^{(k)} + \bar{H}_{\downarrow} \mathbb{I};$$

by setting $\bar{\mathbf{h}}^{(k)}(i_0) = 0$ for some fixed state $i_0 \in \mathcal{S}$. Compute $\tilde{\mathbf{G}}_{\downarrow}^{(k)} = \left(\mathbf{I} - \mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) \right)$. Then $\tilde{\mathbf{G}}_{\downarrow}^{(k)}$ is the matrix with the column corresponding to state i_0 replaced by a column of 1's. Solve the linear system

$$\mathbf{r}_{\downarrow}^{\bar{\boldsymbol{\phi}}^{(k)}} = \tilde{\mathbf{G}}_{\downarrow}^{(k)} \mathbf{w}$$

where $\bar{H}_{\downarrow}^{(k)}$ is the i_0 th component of the solution vector \mathbf{w} and $\bar{\mathbf{h}}^{(k)}(i) = \mathbf{w}(i)$ for $i \neq i_0$;

- 3: **for** $i \in \mathcal{S}$; **do**
- 4: **for** $a \in \mathcal{A}$; **do**
- 5: $f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}^{(k)}) = \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \bar{\mathbf{h}}^{(k)})$;
- 6: Choose $\bar{\boldsymbol{\phi}}^{(k+1)}(i)$ to satisfy

$$\bar{\boldsymbol{\phi}}^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\downarrow}^a(i) + f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}^{(k)}) \right)$$

keeping $\bar{\boldsymbol{\phi}}^{(k+1)}(i) = \bar{\boldsymbol{\phi}}^{(k)}(i)$ when possible;

- 7: **if** $\bar{\boldsymbol{\phi}}^{(k+1)} = \bar{\boldsymbol{\phi}}^{(k)}$ **then**
- 8: Set $\bar{\boldsymbol{\phi}}_{\downarrow}^* = \bar{\boldsymbol{\phi}}^{(k+1)}$ and terminate. Otherwise set $k = k + 1$ and go to Step 2;
- 9: **return** An optimal policy $\bar{\boldsymbol{\phi}}_{\downarrow}^*$;

Algorithm 18 Policy iteration 1 for discrete-time BMDPs with average reward criterion when optimal *upper* bound should be computed

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$;

- 1: Specify $\bar{\phi}^{(1)} \in \Pi$ some pure initial policy, $\bar{\mathbf{h}}^{(0)} = \mathbf{r}_{\downarrow}^{\bar{\phi}^{(1)}}$ and set $k = 1$;
- 2: $\mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) = \arg \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\bar{\phi}^{(k)}}} (\mathbf{P} \bar{\mathbf{h}}^{(k-1)})$;

Solve

$$\mathbf{r}_{\uparrow}^{\bar{\phi}^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) \right) \bar{\mathbf{h}}^{(k)} + \bar{H}_{\uparrow} \mathbf{1};$$

by setting $\bar{\mathbf{h}}^{(k)}(i_0) = 0$ for some fixed state $i_0 \in \mathcal{S}$. Compute $\tilde{\mathbf{G}}_{\uparrow}^{(k)} = \left(\mathbf{I} - \mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(k-1)}) \right)$. Then $\tilde{\mathbf{G}}_{\uparrow}^{(k)}$ is the matrix with the column corresponding to state i_0 replaced by a column of 1's. Solve the linear system

$$\mathbf{r}_{\uparrow}^{\bar{\phi}^{(k)}} = \tilde{\mathbf{G}}_{\uparrow}^{(k)} \mathbf{w}$$

where $\bar{H}_{\uparrow}^{(k)}$ is the i_0 th component of the solution vector \mathbf{w} and $\bar{\mathbf{h}}^{(k)}(i) = \mathbf{w}(i)$ for $i \neq i_0$;

- 3: **for** $i \in \mathcal{S}$; **do**
- 4: **for** $a \in \mathcal{A}$; **do**
- 5: $f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}^{(k)}) = \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \bar{\mathbf{h}}^{(k)})$;
- 6: Choose $\bar{\phi}^{(k+1)}(i)$ to satisfy

$$\bar{\phi}^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\uparrow}^a(i) + f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}^{(k)}) \right)$$

keeping $\bar{\phi}^{(k+1)}(i) = \bar{\phi}^{(k)}(i)$ when possible;

- 7: **if** $\bar{\phi}^{(k+1)} = \bar{\phi}^{(k)}$ **then**
 - 8: Set $\bar{\phi}_{\uparrow}^* = \bar{\phi}^{(k+1)}$ and terminate. Otherwise set $k = k + 1$ and go to Step 2;
 - 9: **return** An optimal policy $\bar{\phi}_{\uparrow}^*$;
-

Algorithm 19 Policy iteration 2 for discrete-time BMDPs with average reward criterion when optimal *lower* bound should be computed

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$;

1: Specify $\bar{\phi}^{(1)} \in \Pi$ some pure initial policy, $\bar{\mathbf{h}}^{(0)} = \mathbf{r}_{\downarrow}^{\bar{\phi}^{(1)}}$ and set $k = 1, l = 1$;

2: **repeat**

3: $\mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) = \arg \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\bar{\phi}^{(k)}}} (\mathbf{P}\bar{\mathbf{h}}^{(l-1)})$;

Solve

$$\mathbf{r}_{\downarrow}^{\bar{\phi}^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) \right) \bar{\mathbf{h}}^{(l)} + \bar{H}_{\downarrow} \mathbf{1};$$

by setting $\bar{\mathbf{h}}^{(l)}(i_0) = 0$ for some fixed state $i_0 \in \mathcal{S}$. Compute

$\tilde{\mathbf{G}}_{\downarrow}^{(k)} = \left(\mathbf{I} - \mathbf{M}_{\downarrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) \right)$. Then $\tilde{\mathbf{G}}_{\downarrow}^{(k)}$ is the matrix with the column corresponding to state i_0 replaced by a column of 1's. Solve the linear system

$$\mathbf{r}_{\downarrow}^{\bar{\phi}^{(k)}} = \tilde{\mathbf{G}}_{\downarrow}^{(k)} \mathbf{w}$$

where $\bar{H}_{\downarrow}^{(k)}$ is the i_0 th component of the solution vector \mathbf{w} and $\bar{\mathbf{h}}^{(l)}(i) = \mathbf{w}(i)$ for $i \neq i_0$;

4: **until** $\|\bar{\mathbf{h}}^{(l)} - \bar{\mathbf{h}}^{(l-1)}\| < \epsilon$;

5: $\bar{\mathbf{h}} = \bar{\mathbf{h}}^{(l)}$;

6: **for** $i \in \mathcal{S}$; **do**

7: **for** $a \in \mathcal{A}$; **do**

8: $f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}) = \min_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \bar{\mathbf{h}})$;

9: Choose $\bar{\phi}^{(k+1)}(i)$ to satisfy

$$\bar{\phi}^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\downarrow}^a(i) + f_{\downarrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}) \right)$$

keeping $\bar{\phi}^{(k+1)}(i) = \bar{\phi}^{(k)}(i)$ when possible;

10: **if** $\bar{\phi}^{(k+1)} = \bar{\phi}^{(k)}$ **then**

11: Set $\bar{\phi}_{\downarrow}^* = \bar{\phi}^{(k+1)}$ and terminate. Otherwise set $k = k + 1, l = 1, \bar{\mathbf{h}}^{(0)} = \bar{\mathbf{h}}$ and go to Step 2;

12: **return** An optimal policy $\bar{\phi}_{\downarrow}^*$;

Algorithm 20 Policy iteration 2 for discrete-time BMDPs with average reward criterion when optimal *upper* bound should be computed

Require: BMDP $(\mathcal{S}, \mathcal{A}, (\mathbf{P}_{\downarrow}^a)_{a \in \mathcal{A}}, (\mathbf{r}_{\downarrow}^a)_{a \in \mathcal{A}})$;

- 1: Specify $\bar{\phi}^{(1)} \in \Pi$ some pure initial policy, $\bar{\mathbf{h}}^{(0)} = \mathbf{r}_{\downarrow}^{\bar{\phi}^{(1)}}$ and set $k = 1, l = 1$;
- 2: **repeat**
- 3: $\mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) = \arg \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^{\bar{\phi}^{(k)}}} (\mathbf{P} \bar{\mathbf{h}}^{(l-1)})$;
Solve

$$\mathbf{r}_{\downarrow}^{\bar{\phi}^{(k)}} = \left(\mathbf{I} - \mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) \right) \bar{\mathbf{h}}^{(l)} + \bar{H}_{\uparrow} \mathbf{1};$$

by setting $\bar{\mathbf{h}}^l(i_0) = 0$ for some fixed state $i_0 \in \mathcal{S}$. Compute $\tilde{\mathbf{G}}_{\uparrow}^{(k)} = \left(\mathbf{I} - \mathbf{M}_{\uparrow}(\mathbf{P}_{\downarrow}^{\bar{\phi}^{(k)}}, \bar{\mathbf{h}}^{(l-1)}) \right)$. Then $\tilde{\mathbf{G}}_{\uparrow}^{(k)}$ is the matrix with the column corresponding to state i_0 replaced by a column of 1's. Solve the linear system

$$\mathbf{r}_{\downarrow}^{\bar{\phi}^{(k)}} = \tilde{\mathbf{G}}_{\uparrow}^{(k)} \mathbf{w}$$

where $\bar{H}_{\uparrow}^{(k)}$ is the i_0 th component of the solution vector \mathbf{w} and $\bar{\mathbf{h}}^{(l)}(i) = \mathbf{w}(i)$ for $i \neq i_0$;

- 4: **until** $\|\bar{\mathbf{h}}^{(l)} - \bar{\mathbf{h}}^{(l-1)}\| < \epsilon$;
- 5: $\bar{\mathbf{h}} = \bar{\mathbf{h}}^{(l)}$;
- 6: **for** $i \in \mathcal{S}$; **do**
- 7: **for** $a \in \mathcal{A}$; **do**
- 8: $f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}) = \max_{\mathbf{P} \in \mathcal{P}_{\downarrow}^a} (\mathbf{P}(i \bullet) \bar{\mathbf{h}})$;
- 9: Choose $\bar{\phi}^{(k+1)}(i)$ to satisfy

$$\bar{\phi}^{(k+1)}(i) = \arg \max_{a \in \mathcal{A}} \left(\mathbf{r}_{\downarrow}^a(i) + f_{\uparrow}^a(\mathbf{P}_{\downarrow}^a(i \bullet), \bar{\mathbf{h}}) \right)$$

keeping $\bar{\phi}^{(k+1)}(i) = \bar{\phi}^{(k)}(i)$ when possible;

- 10: **if** $\bar{\phi}^{(k+1)} = \bar{\phi}^{(k)}$ **then**
 - 11: Set $\bar{\phi}_{\uparrow}^* = \bar{\phi}^{(k+1)}$ and terminate. Otherwise set $k = k + 1, l = 1, \bar{\mathbf{h}}^{(0)} = \bar{\mathbf{h}}$ and go to Step 2;
 - 12: **return** An optimal policy $\bar{\phi}_{\uparrow}^*$;
-

1.4 Solution methods for BSMDPs

Algorithm 21 Transformation method for BSMDPs with average reward criterion

Require: BSMDP $\left((\mathbf{P}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\uparrow}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\uparrow}^a) \right)$;

- 1: **for** $(\rightarrow, \leftarrow) \in \{(\uparrow, \downarrow), (\downarrow, \uparrow)\}$ **do**
- 2: **for** $(i, a) \in \mathcal{S} \times \mathcal{A}$ **do**
- 3: Let the $\{(\mathbf{p}_{\downarrow}^{(i)}, \mathbf{D}_{0\downarrow}^{(i)})\}_{i \in \mathcal{S}}$ and $\{(\mathbf{p}_{\uparrow}^{(i)}, \mathbf{D}_{0\uparrow}^{(i)})\}_{i \in \mathcal{S}}$ be the sets of Phase-type distributions corresponding to the action a ;
- 4: $\mathbf{y}_{\rightarrow}^a(i) = -\mathbf{p}_{\rightarrow}^{(i)} (\mathbf{D}_{0\rightarrow}^{(i)})^{-1} \mathbf{1}$;
- 5: $\mathbf{s}_{\rightarrow}^a(i) = \frac{\mathbf{r}_{\rightarrow}^a(i)}{\mathbf{y}_{\rightarrow}^a(i)}$;
- 6: $\eta = \min_{i \in \mathcal{S}, a \in \mathcal{A}} \frac{\mathbf{y}_{\rightarrow}^a(i)_{\downarrow}}{1 - \mathbf{P}_{\rightarrow}^a(i, i)}$;
- 7: **for** $a \in \mathcal{A}$ **do**
- 8: $\mathbf{S} = \eta (\mathbf{P}_{\rightarrow}^a - \mathbf{I})$;
- 9: $\mathbf{d} = \text{diag}(\mathbf{S})$;
- 10: $\mathbf{R} = \text{diag}(\mathbf{y}_{\rightarrow}^a)^{-1} (\mathbf{S} - \text{diag}(\mathbf{d}))$;
- 11: $\mathbf{e} = \text{diag}(\mathbf{y}_{\leftarrow}^a)^{-1} \mathbf{d} + \mathbf{1}$;
- 12: $\mathbf{Q}_{\rightarrow}^a = \text{diag}(\mathbf{e}) + \mathbf{R}$;
- 13: **return** BMDP $\left((\mathbf{Q}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\uparrow}^a)_{a \in \mathcal{A}} \right), \eta$;

Algorithm 22 Analysis of BSMDPs under the average reward criterion

Require: BSMDP $\left((\mathbf{P}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\uparrow}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\uparrow}^a) \right)$;

- 1: Transform $\left((\mathbf{P}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\uparrow}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\uparrow}^a) \right)$ into BMDP $\left((\mathbf{Q}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\uparrow}^a)_{a \in \mathcal{A}} \right), \eta \in \mathbb{R}$ using Algorithm 21.
- 2: Analyze BMDP $\left(\mathcal{S}, \mathcal{A}, (\mathbf{Q}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\uparrow}^a)_{a \in \mathcal{A}} \right)$ with one of methods 16, 17, 18, 19, 20 and obtain policy ϕ , gain H , bias vector \mathbf{h} .
- 3: **return** $(\phi, H, \eta \mathbf{h})$

Algorithm 23 Transformation method for BSMDPs with discounted reward criterion

Require: BSMDP $\left((\mathbf{P}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\uparrow}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\uparrow}^a) \right)$, discount rate β ;

1: $(\mathbf{Q}_{\uparrow}^a)_{a \in \mathcal{A}} = (\mathbf{P}_{\uparrow}^a)_{a \in \mathcal{A}}$;

2: **for** $a \in \mathcal{A}$ **do**

3: Let the $\{(\mathbf{p}_{\downarrow}^{(i)}, \mathbf{D}_{0\downarrow}^{(i)})\}_{i \in \mathcal{S}}^a$ and $\{(\mathbf{p}_{\uparrow}^{(i)}, \mathbf{D}_{0\uparrow}^{(i)})\}_{i \in \mathcal{S}}^a$ be the sets of Phase-type distributions corresponding to the action a ;

4: **for** $a \in \mathcal{A}$ **do**

5: **for** $i \in \mathcal{S}$ **do**

6: Using $\{(\mathbf{p}_{\downarrow}^{(i)}, \mathbf{D}_{0\downarrow}^{(i)})\}_{i \in \mathcal{S}}^a$ compute lower bound vectors and matrices

$$\mathbf{d}_{1\downarrow}^{(i)} = -\mathbf{D}_{0\downarrow}^{(i)} \mathbf{1};$$

$$\text{Set } \mathbf{P}_{\downarrow}^{(i)} = \mathbf{D}_{0\downarrow}^{(i)} - \beta \mathbf{I} \text{ and } \lambda_{\downarrow} = \max_{\forall i, j \in \mathcal{S}} |\mathbf{P}_{\downarrow}^{(i)}(i, j)|;$$

$$\mathbf{P}_{\downarrow}^{(i)} = \frac{1}{\lambda_{\downarrow}} \mathbf{P}_{\downarrow}^{(i)} + \mathbf{I},$$

$$\mathbf{d}_{1\downarrow}^{(i)} = \frac{1}{\lambda_{\downarrow}} \mathbf{d}_{1\downarrow}^{(i)} \text{ and the time step } \Delta_{\downarrow} = 10/\lambda_{\downarrow};$$

$$\text{Compute } \mathbf{s}_{\downarrow}^a(i) = \mathbf{r}_{\downarrow}^a(i) \int_0^{\infty} (1 - F_{\downarrow}^a(i, t)) e^{-\beta t} dt \text{ and}$$

$$\text{the discount factor } \gamma_{\downarrow}^a(i) = \int_0^{\infty} f_{\downarrow}^a(i, t) e^{-\beta t} dt \text{ with the uniformization based method [1] using } \mathbf{P}_{\downarrow}^{(i)}, \mathbf{d}_{1\downarrow}^{(i)}, \Delta_{\downarrow}.$$

Using $\{(\mathbf{p}_{\uparrow}^{(i)}, \mathbf{D}_{0\uparrow}^{(i)})\}_{i \in \mathcal{S}}^a$ compute upper bound vectors and matrices

$$\mathbf{d}_{1\uparrow}^{(i)} = -\mathbf{D}_{0\uparrow}^{(i)} \mathbf{1};$$

$$\mathbf{P}_{\uparrow}^{(i)} = \mathbf{D}_{0\uparrow}^{(i)} - \beta \mathbf{I};$$

$$\lambda_{\uparrow} = \max_{\forall i, j \in \mathcal{S}} |\mathbf{P}_{\uparrow}^{(i)}(i, j)|;$$

$$\mathbf{P}_{\uparrow}^{(i)} = \frac{1}{\lambda_{\uparrow}} \mathbf{P}_{\uparrow}^{(i)} + \mathbf{I};$$

$$\mathbf{d}_{1\uparrow}^{(i)} = \frac{1}{\lambda_{\uparrow}} \mathbf{d}_{1\uparrow}^{(i)} \text{ and the time step } \Delta_{\uparrow} = 10/\lambda_{\uparrow};$$

$$\text{Compute } \mathbf{s}_{\uparrow}^a(i) = \mathbf{r}_{\uparrow}^a(i) \int_0^{\infty} (1 - F_{\uparrow}^a(i, t)) e^{-\beta t} dt \text{ and}$$

$$\text{the discount factor } \gamma_{\uparrow}^a(i) = \int_0^{\infty} f_{\uparrow}^a(i, t) e^{-\beta t} dt \text{ with the uniformization based method [1] using } \mathbf{P}_{\uparrow}^{(i)}, \mathbf{d}_{1\uparrow}^{(i)}, \Delta_{\uparrow}.$$

7: **return** $(\mathbf{Q}_{\uparrow}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\uparrow}^a)_{a \in \mathcal{A}}, (\gamma_{\uparrow}^a)_{a \in \mathcal{A}}$;

Algorithm 24 Analysis of BSMDPs under the discounted reward criterion

Require: BSMDP $\left((\mathbf{P}_{\dagger}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\dagger}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\dagger}^a) \right)$, discount rate $\gamma \in [0, 1)$;

- 1: Transform $\left((\mathbf{P}_{\dagger}^a)_{a \in \mathcal{A}}, (\mathbf{F}_{\dagger}^a(i, t))_{a \in \mathcal{A}}, (\mathbf{r}_{\dagger}^a) \right)$ into BMDP $\left((\mathbf{Q}_{\dagger}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\dagger})_{a \in \mathcal{A}} \right)$, state-dependent discount factor vector $(\gamma_{\dagger}^a \in \mathbb{R}^{|\mathcal{S}|})_{a \in \mathcal{A}}$ using Algorithm 23.
 - 2: Analyze BMDP $\left(\mathcal{S}, \mathcal{A}, (\mathbf{Q}_{\dagger}^a)_{a \in \mathcal{A}}, (\mathbf{s}_{\dagger})_{a \in \mathcal{A}} \right)$ with one of methods 13, 14, 15 under discount factor $(\gamma_{\dagger}^a)_{a \in \mathcal{A}}$ and obtain policy ϕ , gain vector \mathbf{v} .
 - 3: **return** (ϕ, \mathbf{v})
-

1.5 Solution methods for continuous-time processes

Algorithm 25 Uniformization method for CTMDPs

Require: CTMDP $(\mathcal{S}, \mathcal{A}, \{\mathbf{Q}^a\}_{a \in \mathcal{A}}, \{\mathbf{r}^a\}_{a \in \mathcal{A}})$, discount factor β , *discounted* is true for the discounted reward measure and false else.

- 1: $\lambda = \max_{\forall i \in \mathcal{S}, \forall a \in \mathcal{A}} |\mathbf{Q}^a(i, i)|$;
 - 2: **for** $a \in \mathcal{A}$ **do**
 - 3: $\mathbf{P}^a = \mathbf{I} + \frac{1}{\lambda} \mathbf{Q}^a$;
 - 4: **if** *discounted* **then**
 - 5: $\mathbf{z}^a(i) = \frac{\mathbf{r}^a(i)}{\lambda + \beta}, \quad \forall i \in \mathcal{S}$;
 - 6: **else**
 - 7: $\mathbf{z}^a(i) = \frac{\mathbf{r}^a(i)}{\lambda}, \quad \forall i \in \mathcal{S}$;
 - 8: **if** *discounted* **then**
 - 9: $\gamma = \frac{\lambda}{\lambda + \beta}$;
 - 10: **return** Discrete-time MDP $\{\mathcal{S}, \mathcal{A}, \{\mathbf{P}^a\}_{a \in \mathcal{A}}, \{\mathbf{z}^a\}_{a \in \mathcal{A}}\}$, discount factor γ if *discounted* is true;
-

References

1. D. Gross and D. Miller. The randomization technique as a modeling tool and solution procedure for transient Markov processes. *Operations Research*, 32, 1984.