Aggregation Methods for Markov Decision Problems with Perfect and Imperfect State Information

Kim Hammar (khammar1@asu.edu), Yuchao Li (yuchaoli@asu.edu), and Dimitri P. Bertsekas (dimitrib@mit.edu)

Based on

Section 3.6 of "A Course in Reinforcement Learning: 2nd Edition", by D.P. Bertsekas as well as the upcoming work by Y. Li, K. Hammar, and D.P. Bertsekas

## Aggregation is A Form of Problem Simplification



#### The Aggregation Methodology

- Combine groups of similar states into aggregate states.
- I Formulate an aggregate dynamic programming problem based on these states.
- Solve the aggregate problem using some computational method.
- Use the solution to the aggregate problem to compute a cost function approximation for the original problem.



#### Aggregation for Perfect State Information Problems

2 Aggregation for Imperfect State Information Problems

3 Illustrative Example and Computational Experiments

#### Recap of Markov Decision Problems (MDP)

$$(i) \xrightarrow{p_{ij}(u), g(i, u, j)} (j)$$

- State space:  $X = \{1, \ldots, n\}$ , states are denoted by i, j.
- Control constraint set: U(i).
- Probability of transitioning from state *i* to *j* given control *u*:  $p_{ij}(u)$ .
  - Equivalent formulation:  $x_{k+1} = f(x_k, u_k, w_k)$ .
- Cost of transitioning from state *i* to *j* given control *u*: g(i, u, j).
- Cost-to-go from state i: J(i).
- **Discount** factor:  $\alpha$ .

#### Approximation in Value Space

$$\min_{u \in U(i)} \sum_{j=1}^{n} p_{ij}(u) \left( \overbrace{g(i,u,j)}^{\text{First Step "Future"}} \alpha \widetilde{J}(j) \right)$$

• Optimal policy:  $\mu^*$  can be computed via

$$\mu^*(i) \in \arg\min_{u \in U(i)} \sum_{j=1}^n p_{ij}(u) \big(g(i, u, j) + lpha J^*(j)\big), \quad i = 1, \dots, n$$

where  $J^*$  is the optimal cost function satisfying Bellman's equation

- When computing  $J^*$  is intractable, aggregation computes some  $\widetilde{J}$  to approximate  $J^*$
- A suboptimal policy  $\tilde{\mu}$  can be computed online via

$$ilde{\mu}(i_k) \in rg\min_{u \in U(i_k)} \sum_{j=1}^n p_{i_k j}(u) ig(g(i_k, u, j) + lpha ilde{J}(j)ig)$$

upon reaching state  $i_k$  at stage k

- Introduce a subset A of the original states  $1, \ldots, n$ , called representative states.
- We use *i*, *j* to denote original states and *x*, *y* to denote representative states.



- For each state *i* we define aggregation probabilities  $\{\phi_{ix} \mid x \in A\}$ .
- Intuitively,  $\phi_{ix}$  expresses similarity between states *i* and *x*, where  $\phi_{xx} = 1$ .





## Formulating the Aggregate Dynamic Programming Problem

- State space:  $\mathcal{A}$  (the set of representative states).
- Control constraint set: U(i) (the original control constraint set).
- Transition probabilities and costs

$$\hat{p}_{xy}(u) = \sum_{i=1}^{n} p_{xi}(u) \phi_{iy},$$
 for all representative states  $(x, y)$  and controls  $u$ ,

$$\hat{g}(x,u) = \sum_{i=1}^{n} p_{xi}(u)g(x,u,i)$$

for all representative states x and controls u.



## Solving the Aggregate Dynamic Programming Problem

- The aggregate problem can be solved "exactly" using dynamic programming/simulation; see [Ber19, Section 6.3]
- The optimal cost from a representative state x in this problem is denoted by  $r_x^*$ .



## Cost Difference Between the Aggregate and Original Problems

- The aggregate cost function  $r_x^*$  is only defined for representative states  $x \in A$ .
- The optimal cost function  $J^*(i)$  is defined for the entire state space i = 1, ..., n.
- For a representative state x, we generally have  $r_x^* \neq J^*(x)$ .



## Cost Difference Between the Aggregate and Original Problems

- The aggregate cost function  $r_x^*$  is only defined for representative states  $x \in A$ .
- The optimal cost function  $J^*(i)$  is defined for the entire state space i = 1, ..., n.
- For a representative state x, we generally have  $r_x^* \neq J^*(x)$ .



## Using the Aggregate Solution to Approximate the Original Problem

• We obtain an approximate cost function  $\tilde{J}$  for the original problem via interpolation:

$$\tilde{J}(i) = \sum_{x \in \mathcal{A}} \phi_{ix} r_x^*, \qquad i = 1, \dots, n.$$

• Using this cost function, we can obtain a one-step lookahead policy:

$$ilde{\mu}(i) \in \operatorname*{arg\,min}_{u \in U(i)} \sum_{j=1}^{n} p_{ij}(u) \left( g(i, u, j) + lpha ilde{J}(j) 
ight), \qquad i = 1, 2, \dots, n$$



н	lam	mar	et.	al

## Using the Aggregate Solution to Approximate the Original Problem

#### Approximating the Original Problem

• We obtain an approximate cost function  $\hat{J}$  for the original problem via **interpolation**:

$$ilde{J}(j) = \sum_{y \in \mathcal{A}} \phi_{jy} r_y^*, \qquad \qquad j = 1, \dots, n$$

• Using this cost function, we can obtain a one-step lookahead policy:

#### What is the difference between the approximation $\tilde{J}$ and the optimal cost function $J^*$ ?



# Hard Aggregation

- Consider the case where  $\phi_{jx} = 0$  for all representative states x except one.
- Let  $S_x$  denote the set of states that aggregate to the representative state x.
  - i.e., the *footprint* of x, where  $\{1, \ldots, n\} = \bigcup_{x \in \mathcal{A}} S_x$ .



#### Structure of the Cost Function Approximation

- In the case of hard aggregation,  $\tilde{J}(i) = \sum_{x \in \mathcal{A}} \phi_{ix} r_x^* = r_y^*$  for all  $i \in S_y$ .
- Hence,  $\tilde{J}$  is piecewise constant.



## Approximation Error Bound in the Case of Hard Aggregation [TR96]

• Let  $\epsilon$  be the maximum variation of  $J^*$  within a footprint set  $S_x$ , i.e.,

$$\epsilon = \max_{x \in \mathcal{A}} \max_{i,j \in S_x} |J^*(i) - J^*(j)|.$$

- We refer to the difference  $|J^*(i) \tilde{J}(i)|$  as the approximation error.
- This error is bounded as

$$|J^*(i) - \tilde{J}(i)| \leq \frac{\epsilon}{1-lpha}$$
  $i = 1, \ldots, n.$ 

• Takeway: choose the footprint sets so that  $\epsilon$  is small.



#### General Aggregation and Approximation Error

- Introduce a finite set of aggregate states A.
- Each aggregate state  $x \in A$  is associated with a disjoint subset  $I_x \subset \{1, \ldots, n\}$ .
- An aggregation problem can be defined similarly; see Section 3.6.4 of the course book [Ber25] for details



Upcoming Work

We show that similar error bound also hold for general aggregation with soft aggregation probabilities; i.e.,  $\phi_{jx} \neq 0$  for several  $x \in \mathcal{A}$  [LHB25]

## Partially Observed Markov Decision Problems (POMDPs)

- State space  $X = \{1, ..., n\}$ , observation space Z, and control constraint set U(i).
- Each state transition (*i*, *j*) generates a cost *g*(*i*, *u*, *j*);
- and an observation z with probability p(z | j, u).
- Let b(i) denote the conditional probability that the state is *i*, given the history.
- The belief state is defined as  $b = (b(1), b(2), \dots, b(n))$ .
- The belief b is updated using a **belief estimator** F(b, u, z).
- Goal: Find a policy as a function of b that minimizes the cost.



- The belief *b* resides in the belief space *B*, i.e., the n-1 dimensional unit simplex.
- For example, if the states are  $\{0,1\}$ , then  $b \in [0,1]$ .



• We can obtain representative beliefs via uniform discretization of the belief space:

$$\mathcal{A} = \Big\{ b \mid b \in \mathcal{B}, b(i) = \delta_i / \rho, \sum_{i=1}^n \delta_i = \rho, \delta_i \in \{0, \dots, \rho\} \Big\},$$

where  $\rho$  serves as the discretization resolution.



#### • We can implement hard aggregation via the nearest neighbor mapping:

 $\phi_{by} = 1$  if and only if y is the nearest neighbor of b, where  $b \in B$  and  $y \in A$ .



## Example POMDP: Rocksample (4,3)

- Problem: rover exploration on Mars to find "good" rocks with high scientific value.
- $\bullet\,$  There are 3 rocks on a 4  $\times\,$  4 grid. The rover does not know which rocks are good.
- The controls (north, south, east, west) moves the rover (at cost 0.1).
- The control "sampling" determines the rock quality at the rover position (cost 10 for sampling a bad rock and cost -10 for sampling a good rock).
- Control "check-l" applies a sensor to check the quality of rock / (at cost 1).
- Accuracy of the sensor decreases exponentially with Euclidean distance to the rock.
- The rover stops the mission by moving to the right, yielding an exit-cost of -10.



# Approximating Rocksample (4,3) via Representative Aggregation

- The Rocksample (4,3) POMDP has a 127-dimensional belief space.
- We discretize the belief space with three different resolutions:
  - $\rho = 1$  leads to an aggregate problem with 128 representative beliefs.
  - $\rho = 2$  leads to an aggregate problem with 8256 representative beliefs.
  - $\triangleright$   $\rho=3$  leads to an aggregate problem with 357760 representative beliefs.



Rocksample (4, 3)

# Animation Setup

Rocksample (4, 3)



# Animation Setup

Rocksample (4, 3)



# Animation Setup

Rocksample (4, 3)



Will be presented during the talk.

Will be presented during the talk.

Will be presented during the talk.

- The animations show that performance improves with the discretization resolution  $\rho$ .
- This is not surprising. As  $\rho$  increases,  $\epsilon$  decreases.
- However, the computational complexity increases with the resolution  $\rho$ .



# Comparison Between Aggregation and Other POMDP Methods

POMDP	States n	Observations $ Z $	Controls $ U $	Discount factor $\alpha$
RS (4,4)	257	2	9	0.95
RS (5,5)	801	2	10	0.95
RS (5,7)	3201	2	12	0.95
RS (7,8)	12545	2	13	0.95
RS (10,10)	102401	2	15	0.95

Table: POMDPs used for the experimental evaluation.

Method	Aggregation	Point-based	Heuristic search	Policy-based	Exact DP
Our method	1				
IP [CLZ13]					1
PBVI [PGT06]		1			
SARSOP [Ong+10]		1			
POMCP [SV10]			1		
HSVI [SS12]			1		
AdaOPS [Wu+21]			1		
R-DESPOT [Som+13]			1		
POMCPOW [SK18]			1		
PPO [Sch+17]				1	
PPG [Cob+21]				1	

Table: Methods used for the experimental evaluation; all methods are based on approximation schemes except IP, which uses exact dynamic programming.

Hammar et. al	IDS Cornell Talk	May 7, 2025	27 / 31

# Comparison Between Aggregation and Other POMDP Methods

POMDP Method	RS (4,4)	RS (5,5)	RS (5,7)	RS (7,8)	RS (10, 10)
Aggregation	-17.15/2.4	-18.12/125.5	-17.51/189.1	-14.71/202	-11.59/500
IP	N/A	N/A	N/A	N/A	N/A
PBVI	-8.24/300	-9.05/300	N/A	N/A	N/A
SARSOP	- <b>17.92</b> /10 <sup>-2</sup>	-19.24/58.5	N/A	N/A	N/A
POMCP	-8.64/1.6	-8.80/1.6	-9.81/1.6	-9.46/1.6	-8.98/1.6
HSVI	- <b>17.92</b> /10 <sup>-2</sup>	- <b>19.24</b> /6.2	-24.69/721.3	N/A	N/A
PPO	-8.57/300	-8.15/300	-8.76/300	-7.35/300	-4.59/1000
PPG	-8.57/300	-8.24/300	-8.76/300	-7.35/300	-4.41/1000
AdaOPS	-16.95/1.6	-17.39/1.6	-16.14/1.6	- <b>15.99</b> /1.6	- <b>15.29</b> /1.6
R-DESPOT	-12.07/1.6	-12.09/1.6	-12.00/1.6	-13.14/1.6	-10.41/1.6
POMCPOW	-8.60/1.6	-8.47/1.6	-8.26/1.6	-8.14/1.6	-7.88/1.6

Table: Evaluation results on the benchmark POMDPs; the first number in each cell is the total discounted cost; the second is the compute time in minutes (online methods were given 1 second planning time per control); cells with N/A indicate cases where a result could not be obtained for computational reasons. RS(m,I) stands for an instance of Rocksample with an  $m \times m$  grid and I rocks.

Performance can be further enhanced by combining with other RL methods, such as rollout  $[{\rm Ham}{+}25]$ 

# Thank you!

#### References I

- [Ber19] Dimitri P. Bertsekas. *Reinforcement Learning and Optimal Control*. Athena Scientific, 2019.
- [Ber25] Dimitri P. Bertsekas. *A Course in Reinforcement Learning*. 2nd. Athena Scientific, 2025.
- [CLZ13] Anthony R Cassandra, Michael L Littman, and Nevin Lianwen Zhang. "Incremental pruning: A simple, fast, exact method for partially observable Markov decision processes". In: arXiv preprint arXiv:1302.1525 (2013).
- [Cob+21] Karl W Cobbe et al. "Phasic policy gradient". In: International Conference on Machine Learning. PMLR. 2021, pp. 2020–2027.
- [Ham+25] Kim Hammar et al. "Adaptive Network Security Policies via Belief Aggregation and Rollout". In: *arXiv preprint upcoming* (2025).
- [LHB25] Yuchao Li, Kim Hammar, and Dimitri P. Bertsekas. "Feature-Based Belief Aggregation for Partially Observable Markov Decision Problems". In: arXiv preprint upcoming (2025).
- [Ong+10] Sylvie CW Ong et al. "Planning under uncertainty for robotic tasks with mixed observability". In: *The International Journal of Robotics Research* 29.8 (2010), pp. 1053–1068.

## References II

- [PGT06] Joelle Pineau, Geoffrey Gordon, and Sebastian Thrun. "Anytime point-based approximations for large POMDPs". In: *Journal of Artificial Intelligence Research* 27 (2006), pp. 335–380.
- [Sch+17] John Schulman et al. "Proximal policy optimization algorithms". In: arXiv preprint arXiv:1707.06347 (2017).
- [SK18] Zachary N. Sunberg and Mykel J. Kochenderfer. "Online Algorithms for POMDPs with Continuous State, Action, and Observation Spaces". In: Proceedings of the 28th International Conference on Automated Planning and Scheduling (ICAPS). 2018. URL: https://www.aaai.org/ocs/index. php/ICAPS/ICAPS18/paper/viewFile/17734/16986.
- [Som+13] Adhiraj Somani et al. "DESPOT: Online POMDP planning with regularization". In: Advances in neural information processing systems 26 (2013).
- [SS12] Trey Smith and Reid Simmons. "Heuristic search value iteration for POMDPs". In: *arXiv preprint arXiv:1207.4166* (2012).
- [SV10] David Silver and Joel Veness. "Monte-Carlo Planning in Large POMDPs". In: Advances in Neural Information Processing Systems. Vol. 23. 2010.

#### References III

- [TR96] John N. Tsitsiklis and Benjamin van Roy. "Feature-based methods for large scale dynamic programming". In: Machine Learning 22.1 (Mar. 1996), pp. 59–94.
- [Wu+21] Chenyang Wu et al. "Adaptive online packing-guided search for POMDPs". In: Advances in Neural Information Processing Systems 34 (2021), pp. 28419–28430.