

Example Usage of MDPNv1 Notation

Billy Okal

University of Freiburg

1 Introduction

Many Reinforcement learning (RL) research papers contain paragraphs that define Markov decision process (MDP) and related concepts (see Sutton and Barto, 1998, for detailed exposition). These paragraphs take up space that could otherwise be used to present more useful/new content. In this paper we demonstrate a package implementing a recently proposed notation¹ for MDP that can be used as common foundation. Declaring the use this notation using a single sentence can replace several paragraphs of notational specifications in other papers.

2 MDPs using MDPNv1 notation

Include the package using notation options of: `alpha`, `beta`, `kappa`.

```
1 % ...
2 \usepackage[alpha]{mdpn} % Most verbose
3 %\usepackage[beta]{mdpn} % Compressed
4 %\usepackage[kappa]{mdpn} % Most compressed
5 % ...
```

The MDP is then denoted by a tuple, $(\mathcal{S}, \mathcal{A}, P, \mathcal{R}, R, d_0, \gamma)$, where;

1. We use $t \in \mathbb{N}_{\geq 0}$ to denotes the time step, where $\mathbb{N}_{\geq 0}$ denotes the natural numbers *including zero*.
2. \mathcal{S} is the set of possible states that the agent can be in, and is called the *state set*. The state of the environment at time t is a random variable that we denote by S_t . We will typically use s to denote an element of the state set.
3. Similarly, \mathcal{A} is the set of possible actions the agent can perform. The action at time t is denoted by A_t , while a denotes an element of the action set.
4. \mathcal{R} is the set of possible rewards, defined as $\mathcal{R} \subseteq \mathbb{R} \cup \{-\infty, \infty\}$. Additionally, instantaneous reward at time t is R_t . Elements of the reward set are denoted by r while the infimum and supremum are r_{\min} and r_{\max} respectively.

¹<http://arxiv.org/abs/1512.09075>

5. $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \mapsto [0, 1]$ is called the *transition function*. For all $(s, a, s', t) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathbb{N}_{\geq 0}$, let $P(s, a, s') := \Pr(S_{t+1} = s' \mid S_t = s, t = a)$. That is, P characterizes the distribution over states at time $t + 1$ given the state and action at time t . We allow three alternate notations for P .

- (a) **alpha**: $P(s' \mid s, a) := P(s, a, s')$. This form takes approximately the same amount of space, but makes it more clear that P is a conditional distribution over the next state given the current state and action.
- (b) **beta**: $P_s^a(s') := P(s, a, s')$. This notation moves terms into subscripts and superscripts in order to save some space.
- (c) **kappa**: $P_{s,s'}^a := P(s, a, s')$. This final form is particularly useful when space is limited.

Once the author selects one the three notations modes, consistent within each paper is ensured automatically.

- 6. The reward function is denoted by R , and so on with three options provided.
- 7. d_0 is the initial distribution of states defined as $d_0 : \mathcal{S} \mapsto [0, 1]$.
- 8. γ is the discount factor defined as $\gamma \in [0, 1)$.

Further, π is the policy which is defined as $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1] \dots$

3 Acknowledgments

I want to thank Phillip Thomas for fruitful discussions on the naming and contents of this package.

References

R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.