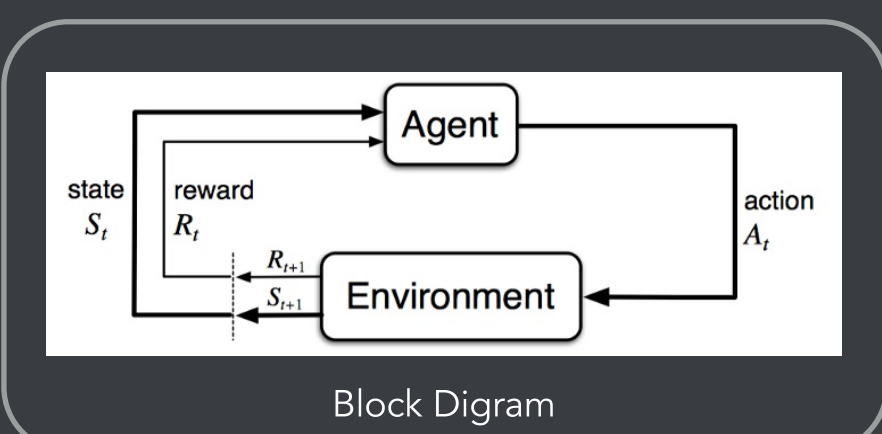
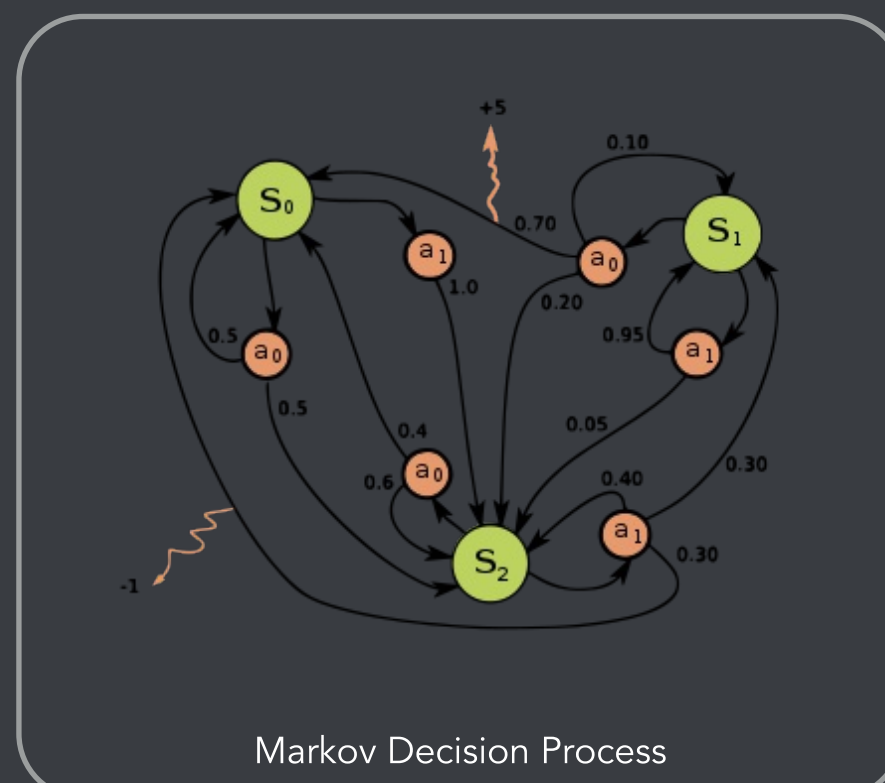


Basic RL Framework



Math Framework



Hypothesis

The Reward Hypothesis: That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).

Reward Hypothesis

Questions

How efficient is to compress all of the useful info for the agent into a single scalar

Important implications

Reward Design

$$\{(S_t, A_t, R_{t+1}, S_{t+1})\}$$

Navigating the MDP Graph produce traces of the form

$$P(S_{t+1}, R_{t+1} | S_t, A_t)$$

MDP Dynamic

Probabilistic Representation is more general as it includes the deterministic one as a special case where distributions are Dirac Delta

 $\pi(A|S)$

Probabilistic Representation is more general as it includes the deterministic one as a special case where distributions are Dirac Delta

Assume you know the First Order MDP Graph exactly which is equivalent to knowing $P(S_{t+1}, R_{t+1} | S_t, A_t)$ then navigating it with a given policy $\pi(A_t | S_t)$ produce traces $\{(S_t, A_t, R_{t+1})_{t=0:T-1}\}$ which means collecting a set of Rewards $\{R_{t+1}\}_{t=0:T-1}$.

Navigating the MDP with a Policy

So the policy π in a state s will get this collection $\{R_t\}_{t \in [0, T]}$

Collecting Rewards

Projecting the Amount of Rewards the Policy is able to collect into the present

Reinforcement Learning in a nutshell

1. Basic RL Framework

1.1. Block Digram

2. Math Framework

2.1. Markov Decision Process

2.1.1. Details

2.1.1.1. Navigating the MDP Graph produce traces of the form

2.1.1.2. MDP Dynamic

2.1.1.2.1. Probabilistic Representation is more general as it includes the deterministic one as a special case where distributions are Dirac Delta

2.1.1.3. Policy

2.1.1.3.1. Probabilistic Representation is more general as it includes the deterministic one as a special case where distributions are Dirac Delta

2.1.1.4. Discounted Reward

2.1.1.4.1. Navigating the MDP with a Policy

2.1.1.4.2. Collecting Rewards

2.1.1.4.3. Projecting the Amount of Rewards the Policy is able to collect into the present

3. Hypothesis

3.1. Reward Hypothesis

3.1.1. Questions

3.1.1.1. How efficient is to compress all of the useful info for the agent into a single scalar

3.1.2. Important implications

3.1.2.1. Reward Design

