

Monte Carlo Learning (MC Learning)

esgl Hu

November 27, 2017

Monte Carlo methods require only *experience* — sample sequences of states, actions, and rewards from actual or simulated interaction with an environment. **It requires no prior knowledge of the environment's dynamics, yet can still attain optimal behavior.**

If a model is not available, then it is particularly useful to estimate *action* values (the values of state-action pairs) rather than *state* values. With a model, state values alone are sufficient to determine a policy; one simply looks ahead one step and chooses whichever action leads to the best combination of reward values and next state. Without a model, however, state values alone are not sufficient.

Algorithm 1 First-visit MC policy evaluation (returns $V \sim v_\pi$)

```
Initialize  $\pi \leftarrow$  policy to be evaluated
Initialize  $V \leftarrow$  an arbitrary state-value function
Initialize  $Returns(s) \leftarrow$  an empty list, for all  $s \in \mathcal{S}$ 
while True do
    Generate an episode using  $\pi$ 
    for each state  $s$  appearing in the episode do
         $G \leftarrow$  return following the first occurrence of  $s$ 
        Append  $G$  to  $Returns(s)$ 
         $V(s) \leftarrow average(Returns(s))$ 
    end for
end while
```
