

Monte Carlo

Alina Vereshchaka

CSE4/510 Reinforcement Learning
Fall 2019

avereshc@buffalo.edu

September 10, 2019

*Slides are based on Deep Reinforcement Learning and Control, CMU 10703, Carnegie-Mellon University

- ▶ **Monte Carlo** methods are learning methods
 - Experience \rightarrow values, policy
- ▶ Monte Carlo uses the simplest possible idea: **value = mean return**
- ▶ Monte Carlo methods can be used in two ways:
 - **Model-free**: No model necessary and still attains optimality
 - **Simulated**: Needs only a simulation, not a full model
- ▶ Monte Carlo methods learn from **complete sample returns**
 - Only defined for episodic tasks (this class)
 - All episodes must terminate (no bootstrapping)

Monte-Carlo Policy Evaluation

- **Goal:** learn $v_{\pi}(s)$ from episodes of experience under policy π

$$S_1, A_1, R_2, \dots, S_k \sim \pi$$

- Remember that the **return** is the total discounted reward:

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

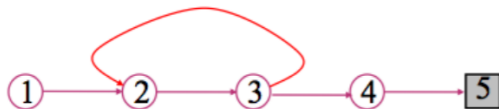
- Remember that the **value function** is the expected return:

$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t \mid S_t = s]$$

- Monte-Carlo policy evaluation uses **empirical mean return** instead of expected return

Monte-Carlo Policy Evaluation

- ▶ **Goal:** learn $v_{\pi}(s)$ from episodes of experience under policy π
- ▶ **Idea:** Average returns observed after visits to s :



- ▶ **Every-Visit MC:** average returns for every time s is visited in an episode
- ▶ **First-visit MC:** average returns only for first time s is visited in an episode
- ▶ Both converge asymptotically

First-Visit MC Policy Evaluation

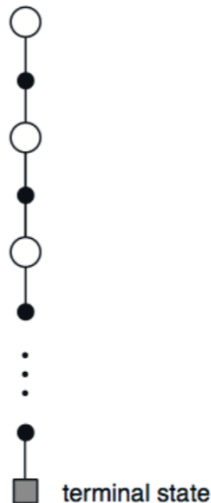
- ▶ To evaluate state s
- ▶ The **first** time-step t that state s is visited in an episode,
- ▶ **Increment counter:** $N(s) \leftarrow N(s) + 1$
- ▶ **Increment total return:** $S(s) \leftarrow S(s) + G_t$
- ▶ Value is estimated by mean return $V(s) = S(s)/N(s)$
- ▶ By law of large numbers $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

Every-Visit MC Policy Evaluation

- To evaluate state s
- **Every** time-step t that state s is visited in an episode,
- **Increment counter:** $N(s) \leftarrow N(s) + 1$
- **Increment total return:** $S(s) \leftarrow S(s) + G_t$
- Value is estimated by mean return $V(s) = S(s)/N(s)$
- By law of large numbers $V(s) \rightarrow v_\pi(s)$ as $N(s) \rightarrow \infty$

Backup Diagram for Monte Carlo

- ▶ Entire rest of episode included
- ▶ Only **one choice** considered at each state (unlike DP)
 - thus, there will be an explore/exploit dilemma
- ▶ Does **not bootstrap** from successor state's values (unlike DP)
- ▶ Value is estimated by **mean return**
- ▶ Time required to estimate one state **does not depend** on the total number of states



Incremental Mean

- ▶ The mean μ_1, μ_2, \dots of a sequence x_1, x_2, \dots can be computed incrementally:

$$\begin{aligned}\mu_k &= \frac{1}{k} \sum_{j=1}^k x_j \\ &= \frac{1}{k} \left(x_k + \sum_{j=1}^{k-1} x_j \right) \\ &= \frac{1}{k} (x_k + (k-1)\mu_{k-1}) \\ &= \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})\end{aligned}$$

Incremental Monte Carlo Updates

- Update $V(s)$ incrementally after episode $S_1, A_1, R_2, \dots, S_T$
- For each state S_t with return G_t

$$N(S_t) \leftarrow N(S_t) + 1$$

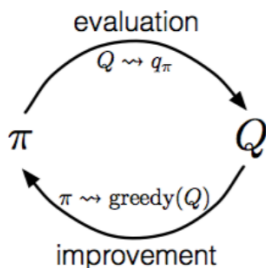
$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

- In non-stationary problems, it can be useful to track a **running mean**, i.e. forget old episodes.

$$V(S_t) \leftarrow V(S_t) + \alpha (G_t - V(S_t))$$

Monte-Carlo Control

$$\pi_0 \xrightarrow{E} q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} q_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_* \xrightarrow{E} q_*$$



- ▶ **MC policy iteration step:** Policy evaluation using MC methods followed by policy improvement
- ▶ **Policy improvement step:** greedify with respect to value (or action-value) function

Monte-Carlo Algorithm

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$Q(s, a) \leftarrow$ arbitrary

$\pi(s) \leftarrow$ arbitrary

$Returns(s, a) \leftarrow$ empty list

Fixed point is optimal
policy π^*

Repeat forever:

Choose $S_0 \in \mathcal{S}$ and $A_0 \in \mathcal{A}(S_0)$ s.t. all pairs have probability > 0

Generate an episode starting from S_0, A_0 , following π

For each pair s, a appearing in the episode:

$G \leftarrow$ return following the first occurrence of s, a

Append G to $Returns(s, a)$

$Q(s, a) \leftarrow \text{average}(Returns(s, a))$

For each s in the episode:

$\pi(s) \leftarrow \arg\max_a Q(s, a)$

Advantages

- MC can be used to learn optimal behavior directly from interaction with the environment. It does not require a model of the environment's dynamics.
- MC can be used with simulation or sample models.
- MC can be used to focus on one region of special interest and be accurately evaluated without having to evaluate the rest of the state set.

Disadvantages

- MC only works for episodic (terminating) environments. It does not work with environment with no terminating states.
- MC must wait until the end of an episode before return is known. For problems with very long episodes this will become too slow.