

Lecture - 12

* First Visit Monte-Carlo (MC)

- intuition
- ① Generate many episodes
 - ② For each state, average the discounted return after the first time it was visited in each episode

Pseudocode:

Initialize:

$\pi \leftarrow$ policy to be evaluated
 $V \leftarrow$ arbitrary state action function

Returns(s) \leftarrow an empty list for all $s \in S$

Repeat forever

- ① Generate an episode using π
- ② For each state s appearing in episode

$G \leftarrow$ Returns following the first occurrence of s

Append G to Returns(s)

$V(s) \leftarrow$ average (Returns(s))

Properties of first visit MC

- Converges almost surely to V^π , if each state is visited infinitely often.

Proof: Consider the seq. of estimates $V_k(s)$ for a particular state s .

$$V_k(s) = \frac{1}{k} \sum_{j=1}^k G_j$$

where, G_j is the j^{th} return from s not the return for j^{th} time

$$\mathbb{E}[G_j] = \mathbb{E}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s, \pi\right]$$

Bias(G_j)

$$= \mathbb{E}[G_j] - \mu$$

$$= V^\pi(s)$$

$$V_1(s) = G_1$$

$$V_2(s) = \frac{G_1 + G_2}{2}$$

\vdots

• Strong law of large numbers.

Let X_1, \dots, X_n be iid random variables with expected value $\mu < \infty$. Then

$\left(\frac{1}{n} \sum_{i=1}^n X_i\right)$ is a sequence of random variables that

converges ~~to~~ almost surely to

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mu$$

$$\therefore P\left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu\right) = 1$$

$\rightarrow \therefore$ By LLN with

$X_i \leftarrow G_i$, we have that

$$V_k(s) \rightarrow V^\pi(s)$$

since, $|s|$ is finite $V_k \rightarrow V^\pi$

How quickly?

$$\text{Var}(V_k(s)) \propto \frac{1}{k}$$

$$\hookrightarrow \text{Var}\left(\frac{1}{k} \sum_{i=1}^k G_i\right)$$

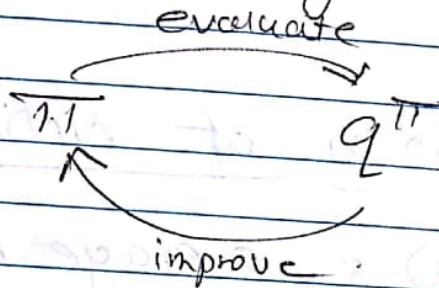
$$= \frac{1}{k^2} \text{Var}\left(\sum_{i=1}^k G_i\right)$$

$$= \frac{1}{k^2} \left[\text{Var}(G_1) + \text{Var}(G_2) + \dots + \text{Var}(G_k) \right]$$

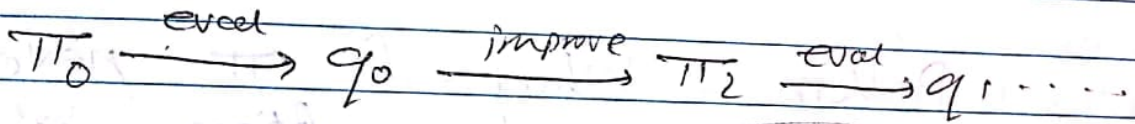
$$= \frac{1}{k^2} \cdot k \text{Var}(G_i) = \frac{1}{k} \text{Var}(G_i)$$

Monte-Carlo Control

- Generalized Policy Iteration (GPI)



① MC policy iteration



- Evaluation using FD-MC
- Use q rather than V .

need to see state-action value infinitely often