

* Model Free Prediction

- Law of Large Number
- Incremental mean
- Monte - Carlo
- Temporal - Difference

* Model Free Prediction

DP에 대한 model의 $R_s^a, P_{ss'}^a$ 를 안지 않음.

$$\begin{aligned} & \rightarrow V_{\pi}(s) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \cdot V_{\pi}(s') \quad \rightarrow V_{k+1}(s) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \cdot V_k(s') \quad \leftarrow \text{Policy eval} \\ & \rightarrow \pi'(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a) \quad \leftarrow \text{Policy improvement (Greedy Policy)} \end{aligned}$$

But generally $R_s^a, P_{ss'}^a$? ... 을라!

→ MC / TD를 써보자 ~

* Law of Large Numbers

= 큰 수의 법칙

R.V X 의 실제 mean $\rightarrow E[X] = \mu \dots$

X_1, X_2, X_3, \dots 가 iid인 X 의 반복적인 독립결과일 때

Also R.V

• Sample mean : $M_n = \frac{1}{n} \sum_{i=1}^n X_i$

• Strong Law of Large Numbers : $P[\lim_{n \rightarrow \infty} M_n = \mu] = 1$

→ 어떤 독립변수들의 Sample mean이든 결국 100% $E[X] = \mu$ 로 수렴!

* Incremental mean

$$\Rightarrow \mu_k = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

$$[PF] \quad \mu_k = \frac{1}{k} \sum_{i=1}^k x_i$$

→ 기존에 저장해 둔 평균값과 새 독립값으로 평균값 갱신!

Coding에 유리!

$$= \frac{1}{k} (x_k + \sum_{i=1}^{k-1} x_i) = \frac{1}{k} (x_k + (k-1) \mu_{k-1})$$

$$\therefore \mu_k = \mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

* Remind

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

$$V_{\pi}(s) = E_{\pi}[G_t | S_t = s]$$

* Monte - Carlo

→ Terminated!

• Episode를 통한 Policy Evaluation

→ Given π 에 따라 episode들로 만들고, 각 episode들을 모두 이용해서 경험적 데이터로 E Get!

Expected return 대신 Empirical mean return!

$$E[G_t] = \sum_k \pi_k R_k(x_k) \quad E[G_t] = \frac{1}{k} \sum_{i=1}^k x_i$$

Law of Large Numbers

• Monte - Carlo Policy Evaluation

Goal : learn V_{π} from episodes

$$S_1, A_1, R_1 \rightarrow S_2, A_2, R_2 \rightarrow \dots$$

① First Visit Monte - Carlo Policy Evaluation

→ 각 Episode에서 내가 관심있는 State 'S'가 나오면,
그 샘플을 기준으로 return을 모두 더해서 mean!

→ 17개의 Episode에 'S'가 여러번 등장? → 첫 등장 시점에서의 G_t 만!

② Every Visit Monte - Carlo Policy Evaluation

→ 각 Episode에서 내가 관심있는 State 'S'가 나오면,
그 샘플을 기준으로 return을 모두 더해서 mean!

→ 17개의 Episode에 'S'가 여러번 등장? → 각 시점에서의 G_t 모두 Sum!

i) Check Episode 1 by 1

ii) if State 's' appeared ... $N(s_t) \leftarrow N(s_t) + 1$
 $S(s_t) \leftarrow S(s_t) + G_t$

iii) Through All episode!

iv) $V(s_t) = \frac{S(s_t)}{N(s_t)}$ With LLN $\rightarrow V(s_t) = V_{\pi}(s_t)$ as $N(s_t) \rightarrow \infty$

⊕ Incremental MC Updates

17개의 Episode로 얻은 $V(s_t)$ 를 Update

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

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t)) \quad (M_k = M_{k-1} + \frac{1}{k} (x_k + M_{k-1}))$$

$$\vdots$$

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

↙ 고정값을 쓰기도 한다.

* Temporal - Difference

→ Not Terminated!

• Episode를 통한 Policy Evaluation

↳ Incomplete Episode를 통한 V_{π} 찾기!

↳ By Bootstrapping

↘ 학습된 모든 info를 원한 값으로 추정!

• A Simple TD Algorithm (TD(0))

Goal : Learn V_{π} online through experience in episode

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$V(s_t) \leftarrow V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

↘ 추정값

TD Target : $R_{t+1} + \gamma V(s_{t+1})$

TD Error : $R_{t+1} + \gamma V(s_{t+1}) - V(s_t) = \delta_t$