# Model-Free Prediction $<\begin{smallmatrix} MC \\ TD \end{smallmatrix}$ (Model Free no knowledge of MDP transitions)

- Basic : MC vs TD vs DP



DP

TD  $V(h_t) \leftarrow V(h_t) + \alpha (R_{t+1} + \gamma V(h_{t+1}) - V(h_t))$

MC  $V(h_t) \leftarrow V(h_t) + \alpha (G_t - V(s_t))$

* Monte-Carlo Policy Evaluation : random sampling → average sample returns,
  - learn $V_\pi$ from episodes of experience under policy $\pi$ (offline)
    ↳ uses empirical mean return ← 끝까지 가서 $G_t$ 구하고 난 후 update ⇒ Sampling의 평균

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

$<$ First-Visit Monte-Carlo  <u>첫음</u> <u>first time-step</u> $t$ that state $S$ is visited in an episode.
  Every-Visit Monte-Carlo  <u>every time-step</u> $t$ that state $s$ is visited in an episode
  아무때나.

To evaluate state $S$

$$N(s) \leftarrow N(s) + 1$$
$$S(s) \leftarrow S(s) + G_t$$

$\Rightarrow V(s) \rightarrow V_\pi(s)$
  as $N(s) \rightarrow \infty$

$$V(s) = S(s) / N(s) \text{ (mean return)}$$

episode by episode, updated

→ MCvsTD        MC                    TD  → MC ⊃ TD라 생각

| MC | TD |
|---|---|
| learn before knowing final outcome | must wait until end of episode before return is known. |
| high variance, zero bias ↳ noisy 하지만 구에 구정. | low variance, some bias ↳ 안정적으로 반비 구정 |
| shallow, sample backup | deep, sample backup. |
| X bootstrap (다계산우 구정) | O bootstrap (구정 구 나아감) |
| O samples | O samples |
|  | n-step에 $n \rightarrow \infty$면 MC로 거걸. |

+) ( Bootstrap . update involves an estimate .
   ( Sample . update samples on expectation

\* Temporal Difference Learning : $n$-step sample → average/weight returns

- learn $v_\pi$ online from experience under policy $\pi$

Value $V(S_t)$

TD error

TD target

$$TD(0) \quad : \quad V(S_t) \leftarrow V(S_t) + \alpha \, ( \, R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \, )$$

$n$-step TD
$$V(S_t) \leftarrow V(S_t) + \alpha \, ( \, G_t^{(n)} - V(S_t) \, )$$

← $n$ step 까지의 return ($n$-return ⇒ $\lambda$-return)

$$( \, G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \quad \cdots \quad + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n}) \, )$$

forward-view TD($\lambda$)
$$V(S_t) \leftarrow V(S_t) + \alpha \, ( \, G_t^{\hat{\lambda}} - V(S_t) \, )$$

$$( \, G_t^{\hat{\lambda}} = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)} \, )$$

Backward-view TD($\lambda$)
$$V(S_t) \leftarrow V(S_t) + \alpha \, \delta_t \, E_t(S)$$

$$( \, \delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t) \; : \; TD \; error \, ) .$$

## * Backward View TD($\lambda$)

→ Eligibility Trace.

$$E_0(s) = 0$$

$$E_t(s) = \underbrace{\gamma\lambda E_{t-1}(s)}_{\text{Frequency}} + \underbrace{1(s_t = s)}_{\text{Recency}}$$

heuristic

→ 정의

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

$$V(s) \leftarrow V(s) + \alpha \underbrace{\delta_t}_{\text{TD error}} \underbrace{E_t(s)}_{\text{eligibility trace}}$$

if (TD(0))

$\lambda = 0$     $E_t(s) = 1(S_t = s)$

(only current
state is updated)     $V(s) \leftarrow V(s) + \alpha \underbrace{\delta_t}_{\text{TD error}} \overset{1}{\cancel{E_t(s)}}$

★ if (TD(1))

$\lambda = 1$     $E_t(s) = \gamma E_{t-1}(s) + 1(s_t = s) = \begin{cases} 0 & \text{if } t < k \\ \gamma^{t-k} & t \geq k \end{cases}$

$$\sum_{t=1}^{T-1} \alpha \delta_t E_t(s) = \alpha \sum_{t=1}^{T-1} \underbrace{\gamma^{t-k}}_{\substack{\text{dis} \\ \text{count}}} \underbrace{\delta_t}_{\text{TD error}} = \alpha(G_k - V(s_k))$$

$$\therefore V(s) \leftarrow V(s) + \alpha(G_k - V(s_k)) \rightarrow \text{Every-Visit MC와 동일!}$$

즉, offline의 MC가 (online) Backward View TD로 환원될수있다.