

< 2. Dynamic Programming Algorithm >

• Contents

2023 2863 배가웅.

- DP 알고리즘 명제
- DP 알고리즘 명제 증명

* DP 알고리즘의 명제

· 초기 state 인 s_1 에 대해서, N -번째부터 1까지 거꾸로 s_t, a_t 를 끌어간다?

↳ $V_1^*(s_1) = V_1(s_1)$ (미래를 알고 있으므로)

<매 t 마다 s_t, a_t 기반해서 비용인 g_t 가 반환됨>

현재 비용도 ↓, 다음 단계 비용도 ↓ ~ 최적의 V

$$V_N(s_N) = g_N(s_N) \quad \text{← } a_N \text{이 존재할 수 없으므로 } s_N \text{만 고려.}$$

$$V_t(s_t) = \min_{a_t \in A_{s_t}} \left[\underbrace{g_t(s_t, a_t)}_{\text{현재 비용}} + \underbrace{V_{t+1}(f(s_t, a_t))}_{a_t \text{함으로써 } t+1 \text{의 비용}} \right] \quad V_t = \{1, \dots, N-1\}$$

· 만약 $a_t^* = d_t^*(s_t)$ 가 $V_t(s_t)$ 식을 최소화하면,

$\pi^* = (d_1^*, d_2^*, \dots, d_{N-1}^*)$ 은 optimal

* DP 알고리즘 명제 증명

· $\pi = (d_1, d_2, \dots, d_{N-1})$ 과 truncated policy $\pi^t = (d_t, d_{t+1}, \dots, d_{N-1})$ 에 대해,

· $V_t^*(s_t)$ 가 최적의 cost-to-go 함수라 가정. → $\pi^* = (d_1^*, \dots, d_{t-1}^*)$ 이 optimal

$$V_t^*(s_t) = \min_{\pi^t} \left\{ g_N(s_N) + \sum_{i=t}^{N-1} g_i(s_i, a_i) \right\}$$

· $t=N$ 일 때, $V_N^*(s_N) = g_N(s_N)$, $V_t^*(s_t) = V_t(s_t)$

별만 방정식에 의해 생성.

· t 와 모든 s_{t+1} 에 대해, $V_{t+1}^*(s_{t+1}) = V_{t+1}(s_{t+1})$ 이라 가정한다면,

$$\begin{aligned} \hookrightarrow \pi^t = (d_t, \pi^{t+1}), \quad \underbrace{V_t^*(s_t)}_{t \text{시점 기준}} &= \min_{(d_t, \pi^{t+1})} \left\{ \underbrace{g_N(s_N)}_{\text{미래}} + \underbrace{g_t(s_t, a_t)}_{\text{현재}} + \underbrace{\sum_{i=t+1}^{N-1} g_i(s_i, a_i)}_{\text{그 사이}} \right\} \\ &= \min_{d_t} \left\{ g_t(s_t, a_t) + \min_{\pi^{t+1}} \left\{ g_N(s_N) + \sum_{i=t+1}^{N-1} g_i(s_i, a_i) \right\} \right\} \\ &\Rightarrow \underbrace{V_{t+1}^*(s_{t+1})}_{\text{미래}} = V_{t+1}(s_{t+1}) \\ &= \min_{d_t} \left\{ g_t(s_t, a_t) + \underbrace{V_{t+1}(s_{t+1})}_{f_t(s_t, a_t)} \right\} \\ &= \min_{a_t \in A_{s_t}} \left\{ g_t(s_t, a_t) + V_{t+1}(f_t(s_t, a_t)) \right\} \\ &= \boxed{V_t(s_t)} \quad \therefore V_1^*(s_1) = V_1(s_1) \quad \square \end{aligned}$$

* State value function, Action value function

· 지금부터 기대되는 Return, 지금 행동을 기대되는 Return
t

* Optimal policy

· 지금부터 기대되는 return을 maximize

* 정의 (가치함수, Optimal policy)

$$\text{Return } G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots$$

$$\textcircled{1} V(s_t) \triangleq \int_{a_t: a_{\infty}} G_t \underbrace{p(a_t, s_{t+1}, a_{t+1}, \dots, | s_t)}_{\text{조건부 확률}} da_t: a_{\infty} \quad (\text{지금부터 시작해서 기대되는 리턴})$$

$$\textcircled{2} \underbrace{Q(s_t, a_t)} \triangleq \int_{s_{t+1}: a_{\infty}} G_t p(s_{t+1}, a_{t+1}, s_{t+2}, a_{t+2}, \dots | s_t, a_t) ds_{t+1}: a_{\infty}$$

$$\textcircled{3} V(s_t) \text{를 maximize 하는 } \begin{pmatrix} p(a_t | s_t) \\ p(a_{t+1} | s_{t+1}) \\ \vdots \\ p(a_{\infty} | s_{\infty}) \end{pmatrix} \text{가 optimal policy.}$$

* Bellman Equation $\textcircled{2}$

$$V(s_t) \triangleq \int_{a_t: a_{\infty}} G_t \underbrace{p(a_t, s_{t+1}, a_{t+1}, \dots | s_t)}_{\textcircled{1} \text{ 베이리안 쿨 적용 } p(x, y) = p(x|y) p(y)}$$

$$p(x, y | z) = p(x | y, z) p(y | z)$$

$$\begin{aligned} \textcircled{1} - \textcircled{1} &= \int_{a_t} \int_{s_{t+1}: a_{\infty}} G_t \underbrace{p(s_{t+1}, a_{t+1}, \dots | s_t, a_t)}_{Q(s_t, a_t)} ds_{t+1}: a_{\infty} p(a_t | s_t) da_t \\ &= \int_{a_t} Q(s_t, a_t) p(a_t | s_t) da_t \end{aligned}$$

$$\textcircled{2} \quad P(a_{t+1}, \dots \mid \cancel{s_t}, \cancel{a_t}, s_{t+1}) \quad p(a_t, s_{t+1} \mid s_t)$$

① - ②