

Bellman Equation

Kihwan Lee

Bellman Equation

Bellman Equation이 무엇인가?

컴퓨터가 계산하기 위해 수학적으로 잘 정리한 것이 MDP

=> 강화학습에서는 이 MDP를 기반

=> 이 MDP를 해결하기 위해 가장 기초적인 개념이 value fn

=> 이 value fn을 계산하기 위한 것이 Bellman Equation

Bellman Equation

Bellman Equation이 무엇인가?

컴퓨터가 계산하기 위해 수학적으로 잘 정리한 것이 MDP

=> 강화학습에서는 이 MDP를 기반

=> 이 MDP를 해결하기 위해 가장 기초적인 개념이 value fn

=> 이 value fn을 계산하기 위한 것이 Bellman Equation

**즉! 강화학습에서 MDP를 해결하기 위해 Bellman Equation을 통해 value fn을 계산하여,
이를 통해 optimal한 policy를 찾는 것이 목표이다.**

Bellman Equation

Bellman Equation이 무엇인가?

컴퓨터가 계산하기 위해 수학적으로 잘 정리한 것이 MDP

=> 강화학습에서는 이 MDP를 기반

=> 이 MDP를 해결하기 위해 가장 기초적인 개념이 value ft

=> 이 value ft을 계산하기 위한 것이 Bellman Equation

State Value Function

$$v_{\pi}(s) = E[G_t | S_t = s]$$

(State) Action Value Function

$$q_{\pi}(s, a) = E[G_t | S_t = s, A_t = a]$$

$$++ \text{Adv ft} = q - v$$

즉! 강화학습에서 MDP를 해결하기 위해 Bellman Equation을 통해 value ft을 계산하여,
이를 통해 optimal한 policy를 찾는 것이 목표이다.

확률 이론 두 가지

1) Law of total probability

$$\bullet P(A) = P(A \cap (\cup B_n)) = \sum_n P(A \cap B_n) = \sum_n P(A|B_n) \cdot P(B_n)$$

$$\bullet E[X] = \sum_x x \cdot P(X=x) = \sum_y E[X|Y=y] \cdot P(Y=y).$$

$$\bullet E[X|Z=z] = \sum_x x \cdot P(X=x|Z=z) = \sum_y E[X|Z=z, Y=y] \cdot P(Y=y|Z=z).$$

2) Law of large numbers

표본의 크기 n 이 커질 수록, 표본 평균이 모평균에 수렴

$$\bar{X}_n = \frac{X_1 + \dots + X_n}{n}, \quad \lim_{n \rightarrow \infty} \bar{X}_n = \mu$$

Bellman Equation 계산

1. Bellman expected Equation

1) State value funtion

$$V_{\pi}(s) = E[G_t | S_t = s] \quad - \text{정의.}$$

$$= \sum_a \underbrace{E[G_t | S_t = s, A_t = a]}_{q_{\pi}(s, a)} \cdot \underbrace{P(A_t = a | S_t = s)}_{\pi(a|s)} \quad - \text{확률 이론 1. (action이 주어졌을 때)}$$

$$= \sum_a \pi(a|s) \cdot q_{\pi}(s, a).$$

$$= \sum_a \pi(a|s) \cdot E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \quad - G_t = R_{t+1} + \gamma G_{t+1} \text{ (return의 정의)}$$

$$= \sum_a \pi(a|s) \cdot \sum_{s', r} E[R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a, S_{t+1} = s', R_{t+1} = r] \cdot P(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a) \quad - \text{확률 이론 1 (s', r이 주어졌을 때)}$$

$$= \sum_a \pi(a|s) \cdot \sum_{s', r} \underbrace{p(s', r | s, a)}_{\text{past}} \cdot [r + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']] \quad - \text{MDP이기에 past는 항상 x.}$$

$$= \sum_a \pi(a|s) \cdot \sum_{s', r} \underbrace{p(s', r | s, a)}_{\text{past}} \cdot [r + \gamma \underbrace{V_{\pi}(s')}]_{\text{present}} \quad - \text{value function 치환.}$$

현재 state s 에서 policy와 function probability가 주어지면 a, r, s' 이다. 즉, 앞의 곱셈은 a, r, s' 이 나온 확률은 곱해준다.

$$= E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s].$$

Bellman Equation 계산

1. Bellman expected Equation

1) State value function

$$V(s_t) = E[G_t | s_t = s]$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]$$

$$= E[R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots) | s_t = s]$$

$$= E[R_{t+1} + \gamma G_{t+1} | s_t = s]$$

$$= E[R_{t+1} + \gamma V(s_{t+1}) | s_t = s]$$

Bellman Equation 계산

1. Bellman expected Equation 을 통해 변경함에 따른 장점?

1) State value function

Bellman expectation Equation은 State-value function을 점화식으로 바꿔주는데 immediate reward R_{t+1} 과 discounted next state value $\gamma v_{\pi}(S_{t+1})$ 의 합으로 분해한 것

=> 이제는 더 이상 return이 필요하지 않습니다! 대신 R_{t+1} 과 S_{t+1} 만 알면 되는 것 입니다.

=> 이는 에피소드가 전체가 다 끝나지 않아도 계산할 수 있다는 장점이 존재

Bellman Equation 계산

1. Bellman expected Equation

2) State action value function

$$q_{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a].$$

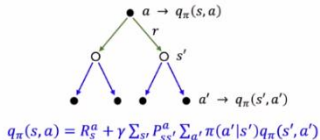
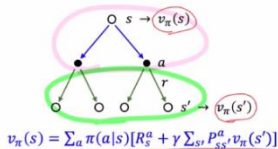
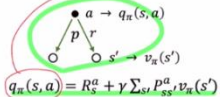
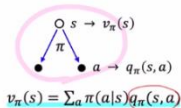
$$= \sum_{s', r} P(s', r | s, a) [r + \gamma \sum_a \pi(a | s') q_{\pi}(s', a)]. \quad - \text{action이 관측될 확률이 policy } \pi.$$

$$= E_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a].$$

Bellman Equation 계산

1. Bellman expected Equation

backup diagrams



Bellman Equation 계산

2. Bellman optimality Equation

value f_t 을 찾는 것도 중요하지만, 우리의 최종 목표는 **reward**를 **최대화** 시키는 **policy** 자체를 찾는 것!

=> 이를 optimal policy라고 하며, 이를 찾기 위해 optimal state value f_t 과 optimal action value f_t 을 이용

=> **bellman optimality equation**

Bellman Equation 계산

2. Bellman optimality Equation

value f_t 을 찾는 것도 중요하지만, 우리의 최종 목표는 **reward를 최대화 시키는 policy** 자체를 찾는 것!

=> 이를 optimal policy라고 하며, 이를 찾기 위해 optimal state value f_t 과 optimal action value f_t 을 이용

=> **bellman optimality equation**

value f_t 들 중에서 maximum이 되는 것이 optimal value f_t . 최적 가치 함수.

이 optimal value f_t 으로 optimal policy를 찾게 되며,

이를 찾게 되는 것이 Markov Decision process를 해결한 것

Bellman Equation 계산

2. Bellman optimality Equation

value f_t 을 찾는 것도 중요하지만, 우리의 최종 목표는 **reward**를 **최대화**시키는 **policy** 자체를 찾는 것!

=> 이를 optimal policy라고 하며, 이를 찾기 위해 optimal state value f_t 과 optimal action value f_t 을 이용

=> **bellman optimality equation**

value f_t 들 중에서 maximum이 되는 것이 optimal value f_t . 최적 가치 함수.

이 optimal value f_t 으로 optimal policy를 찾게 되며,

이를 찾게 되는 것이 Markov Decision process를 해결한 것

=> MDP는 항상 **적어도 하나의 optimal policy가 존재!!**

Bellman Equation 계산

2. Bellman optimality Equation

1) Optimal State value function

$$V_*(s) = \max_{a \in A(s)} q_*(s, a) = \max_a E_{\pi_*} [G_t | S_t = s, A_t = a]$$

$$= \max_a E_{\pi_*} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] \quad - \text{return이 누적되기에 episode 끝날 때까지 고려}$$

$$= \max_a E [R_{t+1} + \gamma V_*(S_{t+1}) | S_t = s, A_t = a] \quad - \text{immediate reward와 next state만 고려}$$

$$\max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')] \quad - \text{sum으로 분해} \quad \Leftrightarrow V_*(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')].$$

immediate reward - r 이 값을
next state s'

Bellman Equation 계산

2. Bellman optimality Equation

1) Optimal state action value function

$$\begin{aligned} q_{\pi}(s, a) &= E_{\pi} [G_t | S_t = s, A_t = a] = E [R_{t+1} + \gamma \max_{a'} q_{\pi}(S_{t+1}, a') | S_t = s, A_t = a]. \\ &= \sum_{s', r} P(s', r | s, a) [r + \gamma \max_{a'} q_{\pi}(s', a')]. \end{aligned}$$

Bellman Equation 계산

2. Bellman optimality Equation

$$v_*(s) = \max_a q_*(s, a)$$

(cf. $v_\pi(s) = \sum_a \pi(a|s) q_\pi(s, a)$)

$$q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a v_*(s')$$

(cf. $q_\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a v_\pi(s')$)

$$v_*(s) = \max_a [R_s^a + \gamma \sum_{s'} P_{ss'}^a v_*(s')]$$

(cf. $v_\pi(s) = \sum_a \pi(a|s) [R_s^a + \gamma \sum_{s'} P_{ss'}^a v_\pi(s')]$)

$$q_*(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \max_{a'} q_*(s', a')$$

(cf. $q_\pi(s, a) = R_s^a + \gamma \sum_{s'} P_{ss'}^a \sum_{a'} \pi(a'|s') q_\pi(s', a')$)