

Real Machine Learning

Объяснение целей через награду

Кумулятивная награда

$$G_t \triangleq R_t + R_{t+1} + R_{t+2} + \dots + R_T$$

Дисконтирование награды

Решение позитивной обратной связи и расходимости

$$G_t \triangleq R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Другая точка зрения: длительность эффекта награды

$$\begin{aligned} G_0 &= R_0 + \gamma R_1 + \gamma^2 R_2 + \dots + \gamma^T R_T \\ &= (1 - \gamma) R_0 \\ &\quad + (1 - \gamma) \gamma (R_0 + R_1) \\ &\quad + (1 - \gamma) \gamma^2 (R_0 + R_1 + R_2) \\ &\quad \dots \\ &\quad + \gamma^T \cdot \sum_{t=0}^T R_t \end{aligned}$$

Решаем вероятность - мат. ожиданием

Наша политика и среда - вероятностные. Максимизируем мат. ожидание!

$$\begin{aligned}\mathbb{E}[G_0] &= \mathbb{E}[R_0 + \gamma R_1 + \dots + \gamma^T R_T] \\&= \mathbb{E}_{E, \pi_\theta}[G_0] \\&= \mathbb{E}_{\pi_\theta}[G_0] \\&= \mathbb{E}[G_0 \mid \pi_\theta] \\&= \mathbb{E}_{\substack{s_0:T \\ a_0:T}}[G_0] \\&= \mathbb{E}_{s_0} \left[\mathbb{E}_{a_0|s_0} \left[R_0 + \mathbb{E}_{s_1|s_0, a_0} \left[\mathbb{E}_{a_1|s_1} [\gamma R_1 + \dots] \right] \right] \right] \\&= \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim p_\theta} [\gamma^t R_t] \\&= \mathbb{E}_{\tau \sim p_\theta(\tau)} [G(\tau)]\end{aligned}$$

$$\mathbb{E} [G_0] = \mathbb{E} [R_0 + \gamma R_1 + \dots + \gamma^T R_T]$$

$$= \mathbb{E}_{E, \pi_\theta} [G_0]$$

$$= \mathbb{E}_{\pi_\theta} [G_0]$$

$$= \mathbb{E} [G_0 \mid \pi_\theta]$$

$$= \mathbb{E}_{\substack{s_0:T \\ a_0:T}} [G_0]$$

$$= \mathbb{E}_{s_0} \left[\mathbb{E}_{a_0|s_0} \left[R_0 + \mathbb{E}_{s_1|s_0, a_0} \left[\mathbb{E}_{a_1|s_1} [\gamma R_1 + \dots] \right] \right] \right]$$

$$= \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim p_\theta} [\gamma^t R_t]$$

$$= \mathbb{E}_{\tau \sim p_\theta(\tau)} [G(\tau)]$$

$$\tau \triangleq (s_0, a_0, s_1, \dots, a_{T-1}, s_T)$$

$$p_\theta(\tau) = p(s_0) \prod_{t=0}^{T-1} \pi_\theta(a_t|s_t) p(s_{t+1}|s_t, a_t)$$

Уравнения Беллмана

Для постановки задачи динамического программирования

$V(s)$ - по политике

$$\begin{aligned} v_{\pi}(s) &= \sum_a \pi(a \mid s) \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma v_{\pi}(s')] \\ &= \mathbb{E}_{\pi} [R_t + \gamma v_{\pi}(S_{t+1}) \mid S_t = s] \end{aligned}$$

$q(s, a)$ - по политике

$$\begin{aligned} q_{\pi}(s, a) &= \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma v_{\pi}(s')] \\ &= \sum_{r, s'} p(r, s' \mid s, a) \left[r + \gamma \sum_{a'} \pi(a' \mid s') q_{\pi}(s', a') \right] \end{aligned}$$

$V(s)$ - оптимальное

$$\begin{aligned} v_*(s) &= \max_a \sum_{r, s'} p(r, s' \mid s, a) [r + \gamma v_*(s')] \\ &= \max_a \mathbb{E} [R_t + \gamma v_*(S_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

$q(s, a)$ - оптимальное

$$\begin{aligned} q_*(s, a) &= \mathbb{E} \left[R_t + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a \right] \\ &= \sum_{r, s'} p(r, s' \mid s, a) \left[r + \gamma \max_{a'} q_*(s', a') \right] \end{aligned}$$

Операторный вид уравнений Беллмана

$$[\mathcal{T}^\pi V](s) = \mathbb{E}_{r,s'|s,a=\pi(s)} [r + \gamma V(s')]$$

$$[\mathcal{T}^\pi Q](s, a) = \mathbb{E}_{r,s'|s,a} [r + \gamma \mathbb{E}_{a' \sim \pi(s')} [Q(s', a')]]$$

$$[\mathcal{T}V](s) = \max_a \mathbb{E}_{r,s'|s,a} [r + \gamma V(s')]$$

$$[\mathcal{T}Q](s, a) = \mathbb{E}_{r,s'|s,a} \left[r + \gamma \max_{a'} Q(s', a') \right]$$

Все операторы порождают сжимающие отображения - значит имеют неподвижную точку

$$[\mathcal{T}^\pi V](s) = \mathbb{E}_{r,s'|s,a=\pi(s)} [r + \gamma V(s')]$$

$$[\mathcal{T}^\pi Q](s, a) = \mathbb{E}_{r,s'|s,a} [r + \gamma \mathbb{E}_{a' \sim \pi(s')} [Q(s', a')]]$$

$$[\mathcal{T}V](s) = \max_a \mathbb{E}_{r,s'|s,a} [r + \gamma V(s')]$$

$$[\mathcal{T}Q](s, a) = \mathbb{E}_{r,s'|s,a} \left[r + \gamma \max_{a'} Q(s', a') \right]$$

Не можете оценить – не можете улучшить

Оценка качества политики

$$\begin{aligned} v_{\pi}(s) &= \sum_a \pi(a | s) \sum_{r, s'} p(r, s' | s, a) [r + \gamma v_{\pi}(s')] \\ &= \mathbb{E}_{\pi} [R_t + \gamma v_{\pi}(S_{t+1}) | S_t = s] \end{aligned}$$

Улучшение политики

$$\pi'(s) \leftarrow \arg \max_a \overbrace{\sum_{r, s'} p(r, s' | s, a) [r + \gamma v_\pi(s')]}^{q_\pi(s, a)}$$

$$v_{\pi'}(s) = \max_a \sum_{r, s'} p(r, s' | s, a) [r + \gamma v_\pi(s')]$$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Repeat

$\Delta \leftarrow 0$

For each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement

policy-stable \leftarrow *true*

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow *false*

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Initialize array V arbitrarily (e.g., $V(s) = 0$ for all $s \in \mathcal{S}^+$)

Repeat

$$\Delta \leftarrow 0$$

For each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, $\pi \approx \pi_*$, such that

$$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

Model-free

Не знаем

$$P(s', r | s, a)$$

Монте карло

$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} [r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')]$$

Q - можно приближать итеративно

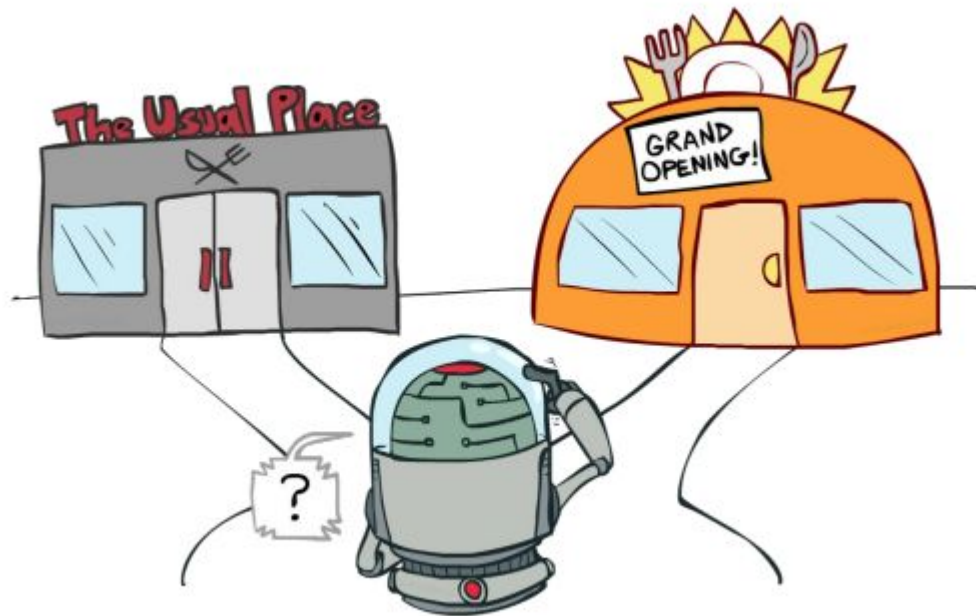
$$Q(s_t, a_t) \leftarrow E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a')$$

$$E_{r_t, s_{t+1}} r_t + \gamma \cdot \max_{a'} Q(s_{t+1}, a') \approx \frac{1}{N} \sum_i r_i + \gamma \cdot \max_{a'} Q(s_i^{next}, a')$$

$$Q(s_t, a_t) \leftarrow \alpha \cdot \hat{Q}(s_t, a_t) + (1 - \alpha) Q(s_t, a_t)$$

$$\pi(s) : \operatorname{argmax}_a Q(s, a)$$

Исследование и эксплуатация



Жадно и не очень

Жадно – эpsilon жадно!

$$\pi(a|s) = \textit{softmax}(\frac{Q(s,a)}{\tau})$$

on-policy & off-policy

REINFORCE

Expected reward: $R(z)$ setting

$$J = E_{\substack{s \sim p(s) \\ a \sim \pi_{\theta}(s|a) \\ \dots}} R(s, a, s', a', \dots)$$

Expected discounted reward: $R(s,a) = r + \gamma * R(s',a')$

$$J = E_{\substack{s \sim p(s) \\ a \sim \pi_{\theta}(s|a)}} Q(s, a)$$

“true” Q-function

$$J = E_{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}} Q(s, a) = \int_s p(s) \int_a \pi_\theta(a|s) Q(s, a) da ds$$

$$J = \underset{\substack{s \sim p(s) \\ a \sim \pi_{\theta}(s|a)}}{E} Q(s, a) = \int_s p(s) \int_a \pi_{\theta}(a|s) Q(s, a) da ds$$

$$J \approx \frac{1}{N} \sum_{i=0}^N \sum_{s, a \in z_i} Q(s, a)$$


True action value
a.k.a. $E[R(s, a)]$

sample N sessions

$$\nabla \log \pi(z) = \frac{1}{\pi(z)} \cdot \nabla \pi(z)$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$

$$\nabla J = \int_s p(s) \int_a \nabla \pi_\theta(a|s) Q(s, a) da ds$$

$$\pi \cdot \nabla \log \pi(z) = \nabla \pi(z)$$


$$\nabla J = E_{\substack{s \sim p(s) \\ a \sim \pi_\theta(s|a)}} \nabla \log \pi_\theta(a|s) \cdot Q(s, a)$$