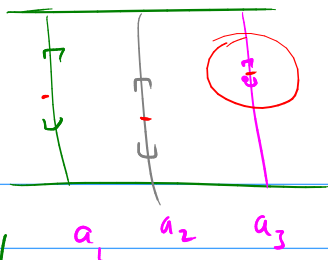
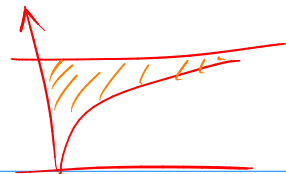


$a_1 \dots a_n$



regret



Contextual bandits

+ context $c \in C$

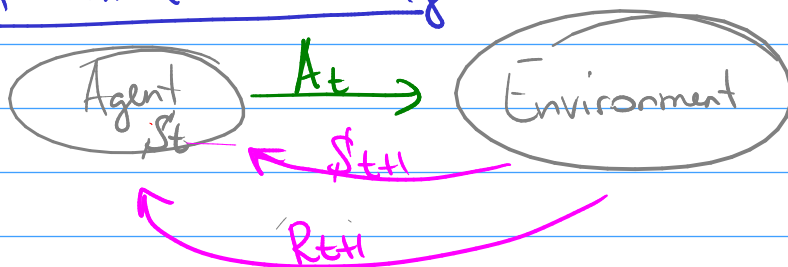
A/B testing

$$E_p[f] = \int p(x) f(x) dx = \int q(x) \left[\frac{p(x)}{q(x)} f(x) \right] dx =$$

$$= E_q \left[f \cdot \frac{p}{q} \right]$$

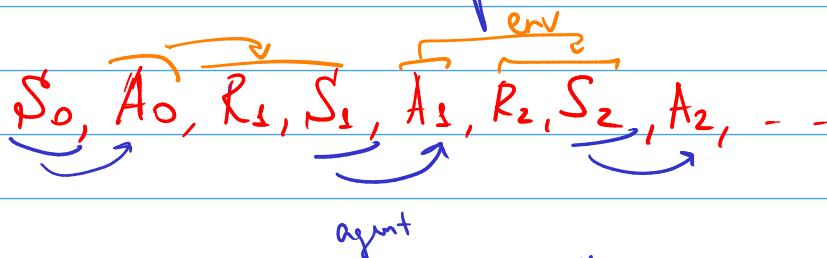
importance weights

Reinforcement Learning



MDP - Markov decision process

Trajectory:



S, R
 $A(s)$

Strategy
Dynamics

$$\pi(a|s)$$

$$\pi: S \rightarrow \text{Prob}[A(s)]$$

$$p(s', r | s, a) = p(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$$

$$p(s' | s, a)$$

Return

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

$$E[G_t] \rightarrow \max$$

$$R_{t+1}, R_{t+2}, \dots, R \dots$$

$\gamma < 1$

discount

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[\sum_k \gamma^k R_{t+k+1} | S_t = s \right]$$

$$Q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[\sum_k \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

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

$$V_{\pi}(s) = \mathbb{E}_{a \sim \pi(s)} [Q_{\pi}(s, a)]$$

$$Q_{\pi}(s, a) = r(s, a) + \gamma \cdot \mathbb{E}_{s' \sim p(s'|s, a)} [V_{\pi}(s')]]$$

Bellman equations

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] =$$

$$= \mathbb{E}_{a \sim \pi(a|s)} \left[R_{t+1} + \gamma \mathbb{E}_{\pi} [G_{t+1}] | S_t = s \right] =$$

$$= \sum_a \pi(a|s) \sum_{s', z} p(s', z | s, a) \cdot \left(z + \gamma \cdot \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s'] \right)$$

$$\boxed{V_{\pi}(s) = \sum_a \pi(a|s) \sum_{s', z} p(s', z | s, a) \left(z + \gamma \cdot V_{\pi}(s') \right)} \quad \begin{array}{l} \text{Bellman} \\ \text{eq} \end{array}$$

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

$$= \sum_{s', z} p(s', z | s, a) \cdot \left(z + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s'] \right)$$

$$\boxed{Q_{\pi}(s, a) = \sum_{s', z} p(s', z | s, a) \left(z + \gamma \sum_{a'} \pi(a'|s') Q_{\pi}(s', a') \right)} \quad \begin{array}{l} \text{Bellman} \\ \text{eq} \end{array}$$

$$RL? \quad \left(\pi^* = \operatorname{argmax}_{\pi} V_{\pi}(s) \right)$$

$$V_{\pi}(s), s \in \mathcal{S}$$

$$\overline{V_{\pi}(s)} = A \cdot \overline{V_{\pi}(s)} + \overline{b}$$

$$\overline{x} = f(\overline{x})$$

$$\overline{z^{(k)}} = f(\overline{z^{(k-1)}})$$

① Найти велел. π , когда V_x, Q_x

② Не знаем $p(s', z | s, a)$

③ Дать $|S|$ асимптотически, где функция ур. берем по порядку

$$V_x(s) = \max_{\pi} V_{\pi}(s)$$

optimal value function

$$Q_x(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

$$\pi_x(s) = \operatorname{argmax}_a Q_x(s, a) = \operatorname{argmax}_a \sum_{s', z} p(s', z | s, a) (z + \gamma V_x(s'))$$

$$\pi_x^*(s) = \operatorname{argmax}_a Q_x(s, a)$$

PLANNING

$$V_x(s) = \max_{\pi} V_{\pi}(s) = \max_{\pi} E_{\pi} [G_t | S_t = s] = \max_a E [R_{t+1} + \gamma \max_{\pi} V_{\pi}(S_{t+1})]$$

$$V_x(s) = \max_a \sum_{s', z} p(s', z | s, a) (z + \gamma V_x(s'))$$

Bellman optimality equation

$$Q_x(s, a) = \max_{\pi} E [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a] =$$

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

$$\max_{a', \pi} Q_{\pi}(S_{t+1}, a') = \max_{a'} Q_*(S_{t+1}, a')$$

$$Q_*(s, a) = \sum_{s', z} p(s', z | s, a) (r + \gamma \max_{a'} Q_*(s', a'))$$

$Q_2 f(Q)$

Policy improvement

$$\pi \rightsquigarrow \pi'$$

$$\forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

Policy improvement theorem

Esau $\forall s \quad Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s)$
 TO $\forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$

$$\begin{array}{c} \pi \quad \pi'(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a) \\ \downarrow \\ Q_{\pi}(s, a) \end{array}$$

1. bo:

$$V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s)) =$$

$$= \mathbb{E} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = \pi'(s)]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s] \leq$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma Q_{\pi}(S_{t+1}, \pi'(S_{t+1})) | S_t = s] =$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma \mathbb{E} [R_{t+2} + \gamma V_{\pi}(S_{t+2}) | S_{t+1}, A_{t+1} = \pi'(S_{t+1})] | S_t = s]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(S_{t+2}) | S_t = s] \leq \dots$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s] = V_{\pi'}(s)$$

Policy Iteration

$$\pi^{(0)} \rightarrow V^{(0)} \rightarrow \pi^{(1)} \rightarrow V^{(1)} \rightarrow \dots$$

Value Iteration

$$\begin{aligned} V^{(k+1)}(s) &= \sum_a \pi^{(k)}(a|s) \sum_{s', z} p(s', z | s, a) (r + \gamma V^{(k)}(s')) \\ &= \max_a \sum_{s', z} p(s', z | s, a) (r + \gamma V^{(k)}(s')) \end{aligned}$$

$$\pi^{(k+1)}(s) = \operatorname{argmax}_a Q_{\pi^{(k)}}(s, a) =$$

$$= \operatorname{argmax}_a \sum_{s', z} p(s', z | s, a) (z + \gamma V^{(k)}(s'))$$

$$V^{(k+1)}(s) = \max_a \sum_{s', z} p(s', z | s, a) (z + \gamma V^{(k)}(s'))$$

vs $\pi(s) = \operatorname{argmax}_a Q_{\pi}(s, a)$

$$V(s) = \max_a \sum_{s', z} p(s', z | s, a) (z + \gamma V(s'))$$

$$\Rightarrow V = V_{\pi}$$

2

$p(s', z | s, a) = ?$

Traj $s_0, a_0, z_1, s_1, a_1, z_2, s_2, \dots$

Monte Carlo

V, Q overcame by Thompson

input: π

$Q(s, a)$ $Ret(s, a)$

output: $V_{\pi}(s)$

MC prediction

- init $V(s), Returns(s)$

- loop:

- generate episode $s_0, a_0, R_1, s_1, a_1, \dots, R_T$ no π

- compute $G_T = 0, G_t = \gamma G_{t+1} + R_{t+1}$

- append G_t to $Returns(s_t)$

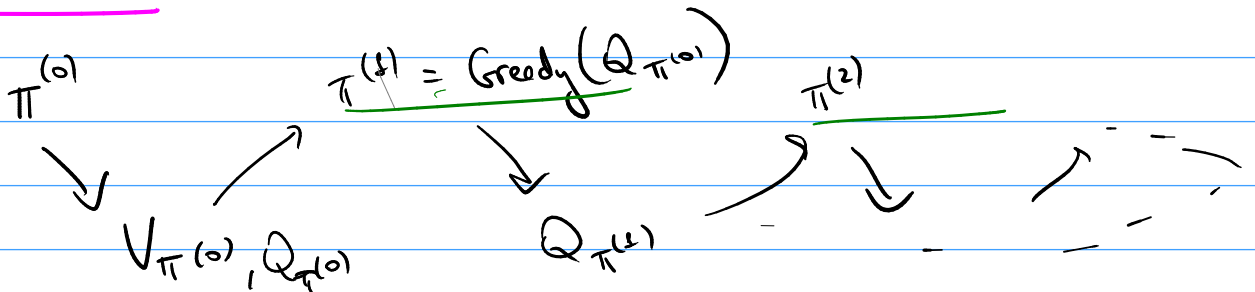
- $V(s_t) = \text{Avg}(Returns(s_t))$

$Ret(s_t, a_t)$

$Ret(s_t, a_t)$

first-visit MC
every-visit MC

Exploring starts



loop:

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

- $Q(s, a) := \text{MCpred}(\pi)$

prediction

$\pi \rightarrow V_{\pi}, Q_{\pi}$

control

π

MC control

- init π, Q, Ret
- loop:
 - gen. ep. by π : $S_0, A_0, R_1, \dots, S_T, A_{T-1}, R_T$
 - compute $G_t, t=0, T-1$
 - $\forall t$: append G_t to $\text{Ret}(S_t, A_t)$
- $Q(S_t, A_t) \leftarrow \text{Avg}(\text{Ret}(S_t, A_t))$
- $\pi(S_t) := \underset{a}{\text{argmax}} Q(S_t, a)$

Soft policy

$$\pi(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A(s)|}, & a = a_* \\ \frac{\epsilon}{|A(s)|}, & \forall a \neq a_* \end{cases}$$

ϵ -greedy

π' - ϵ -greedy

π - ϵ -soft

$$\pi(a|s) \rightarrow \frac{\epsilon}{|A(s)|} \quad \forall a$$

$$Q_{\pi'}(s, \pi'(s)) = \sum_a \pi'(a|s) Q_{\pi}(s, a) =$$

$$= \frac{\epsilon}{|A(s)|} \sum_{a \neq a_*} Q_{\pi}(s, a) + \left(1 - \epsilon + \frac{\epsilon}{|A(s)|}\right) Q_{\pi}(s, a_*) =$$

$= \max_{a'} Q_{\pi}(s, a')$

$$= \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + (1 - \epsilon) Q_{\pi}(s, a_*) =$$

$$= \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + (1 - \epsilon) \max_{a'} Q_{\pi}(s, a')$$

$\sum_a \left(\frac{\pi(a|s) - \frac{\epsilon}{|A(s)|}}{1 - \epsilon}\right) = 1$

$$V_{\pi'}(s) = Q_{\pi'}(s, \pi'(s)) = \frac{\epsilon}{|A(s)|} \sum_a Q_{\pi}(s, a) + \sum_a \left(\frac{\pi(a|s) - \frac{\epsilon}{|A(s)|}}{1 - \epsilon}\right) \cdot Q_{\pi}(s, a)$$

$$(1 - \epsilon) \cdot \sum_a \left(\frac{\pi(a|s) - \frac{\epsilon}{|A(s)|}}{1 - \epsilon}\right) \cdot Q_{\pi}(s, a)$$

$$Q_{\pi}(s, \pi'(s)) \geq V_{\pi'}(s)$$