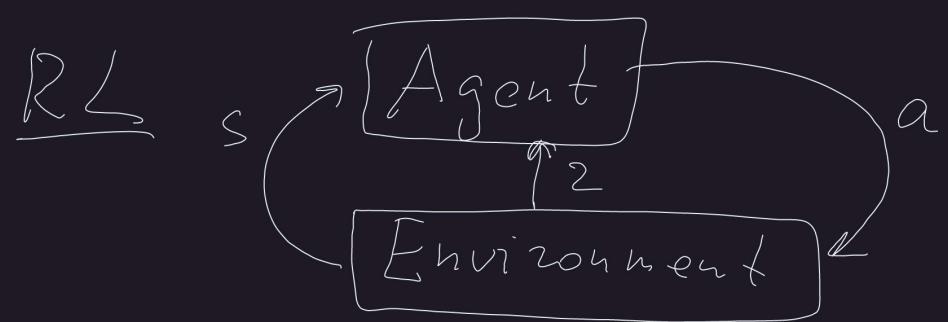


ML/DL NNs for control: Reinforcement Learning

Dataset $\{x_i, y_i\}_{i=1}^N \quad f(x, \theta)$

$$F(\theta) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f(x_i, \theta)) + \lambda R(\theta) \rightarrow \min_{\theta}$$

$$\theta = \text{opt_step}(\nabla_{\theta} F)$$



Trajectory

$$\tau = \{ s_0, a_0, r_0, s_1, a_1, r_1, \dots \}$$

$$p(\tau | \pi) = p(s_0) \pi(a_0 | s_0) p(s_1 | s_0, a_0) \pi(a_1 | s_1) \dots$$

MDP

$$p(s_t | a_{t-1}, s_{t-1}, a_{t-2}, s_{t-2}, \dots, a_0, s_0) = p(s_t | s_{t-1}, a_{t-1})$$

$s \in S$ - state

$a \in A$ - action

$\pi(a|s)$ - policy

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

$r(s, a)$ - reward

$$R_t := \gamma_t + \gamma \gamma_{t+1} + \gamma^2 \gamma_{t+2} + \dots \quad \gamma \in (0, 1)$$

\uparrow
cumulative reward or
return

$$J(\pi) = \underset{P(\pi)}{\mathbb{E}} R_0 \rightarrow \max \pi$$

$$V^\pi(s) := \mathbb{E}_{\substack{p(a|s)}} [R_t \mid s_t = s] \quad - \text{value function}$$

or state function

$$V^*(s) := \max_{\pi} V^\pi(s)$$

$$Q^\pi(s, a) := \mathbb{E}_{\substack{p(a|s)}} [R_t \mid s_t = s, a_t = a] \quad - \text{quality func.}$$

or state-action func.

$$Q^*(s, a) = \max_{\pi} Q^\pi(s, a)$$

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots = r_t + \gamma (r_{t+1} + \gamma r_{t+2} + \dots) =$$

$$= r_t + \gamma R_{t+1}$$

$$V^\pi(s_0) = E[R_0] = E_{a_0, s_1, a_1, \dots} (r(s_0, a_0) + \gamma R_1) =$$

$$= E_{\pi(a_0|s_0)} (r(s_0, a_0) + \gamma E_{P(s_1|s_0, a_0)} (E_{a_1, s_2, a_2, \dots} R_2))$$

$$\Rightarrow V^\pi(s) = E_{\pi(a|s)} (r(s, a) + \gamma E_{P(s'|s, a)} V^\pi(s'))$$

$$V^\pi(s_1)$$

$$\begin{aligned}
V^*(s_0) &= \max_{\pi} V^\pi(s_0) = \\
&= \max_{\pi_0(a_0|s_0), \pi_1(a_1|s_1), \dots} \left[\mathbb{E}_{\pi_0(a_0|s_0)} \left(r(s_0, a_0) + \gamma \mathbb{E}_{p(s_1|s_0, a_0)} V^\pi(s_1) \right) \right] = \\
&= \left\{ \max(a + \beta, a + c) = a + \max(\beta, c) \right\} = c V^*(s_1) \\
&= \max_{\pi_0(a_0|s_0)} \left[\mathbb{E}_{\pi_0(a_0|s_0)} \left(r(s_0, a_0) + \gamma \mathbb{E}_{p(s_1|s_0, a_0)} \left(\max_{\pi_1(a_1|s_1), \pi_2, \dots} V^\pi(s_1) \right) \right) \right] = \\
&= \max_{\pi_0} \mathbb{E}_{\pi_0(a_0|s_0)} \left(\underbrace{r(s_0, a_0) + \gamma \mathbb{E}_{p(s_1|s_0, a_0)} V^*(s_1)}_{= y(s_0, a_0)} \right)
\end{aligned}$$

$$E_{\pi_o(a_o | s_o)} y(s_o, a_o) = \sum_{a_o} \pi_o(a_o | s_o) y(s_o, a_o) \rightarrow \max_{\pi_o}$$

$$\pi^{opt}(a_o | s_o) = \underset{a_o}{\arg \max} \left[y(s_o, a_o) \right] \quad \left\{ \begin{array}{l} \pi_o(a_o | s_o) \geq 0 \quad \forall a_o \\ \sum_{a_o} \pi_o(a_o | s_o) = 1 \end{array} \right.$$

$$V^*(s) = \max_a \left(r(s, a) + \gamma E_{p(s'|s, a)} V^*(s') \right)$$

$$\pi^*(a | s) = \underset{a}{\arg \max} \left(- - - \right)$$

$$Q^{\pi}(s, a) = \gamma(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \mathbb{E}_{\pi(a'|s')} Q^{\pi}(s', a')$$

$$Q^*(s, a) = \gamma(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a')$$

$$V^{\pi}(s) = \mathbb{E}_{\pi(a|s)} Q^{\pi}(s, a)$$

$$Q^{\pi}(s, a) = \gamma(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^{\pi}(s')$$

$$V^*(s) = \max_a Q^*(s, a)$$

$$Q^*(s, a) = \gamma(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^*(s')$$

Value Iteration

$$V^*(s) = \max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^*(s')) \quad \text{for } s$$

$$V^* = F(V^*)$$

$$V_{\text{new}} = F(V_{\text{old}}) \quad - \text{fixed-point iteration}$$

$$\|F(V) - F(W)\| < \|V - W\| \quad \forall V, W$$

Init. V

For iterations:

$$\delta = 0$$

For all s :

$$v = V(s)$$

$$V(s) = \max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V(s'))$$

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

if $\Delta < \varepsilon$, then break

Inference:

$$\pi(a|s) = \arg\max_a (r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V(s'))$$

2 word:

| | | | |
|---|---|----|---|
| 0 | 0 | 0 | 1 |
| 0 | 0 | -1 | |
| 0 | 0 | 0 | 0 |

terminal state

s - position

$$a \in \{\leftarrow, \rightarrow, \uparrow, \downarrow\}$$

V_0 :

| | | | |
|---|---|---|----|
| 0 | 0 | 0 | 1 |
| 0 | 2 | 0 | -1 |
| 0 | 0 | 0 | 0 |

\rightarrow

V_1 :

| | | | |
|---|---|------------|----|
| 0 | 0 | γ | 1 |
| 0 | 2 | γ^2 | -1 |
| 0 | 0 | 0 | 0 |

$\rightarrow V_2$:

| | | | |
|------------|------------|------------|----|
| γ^7 | γ^2 | γ | 1 |
| γ^6 | γ^3 | γ^2 | -1 |
| γ^5 | γ^4 | γ^3 | 0 |

\rightarrow

$\rightarrow V_3$:

| | | | |
|------------|------------|------------|------------|
| γ^3 | γ^2 | γ | 1 |
| γ^6 | γ^3 | γ^2 | -1 |
| γ^5 | γ^4 | γ^3 | γ^4 |

$\rightarrow V_4$:

| | | | |
|------------|------------|------------|------------|
| γ^3 | γ^2 | γ | 1 |
| γ^4 | γ^3 | γ^2 | -1 |
| γ^5 | γ^4 | γ^3 | γ^4 |

$= V_{\text{opt}}$

$$\max(0 + \gamma \cdot 0, 0 + \gamma(-1)) = 0$$

Q-learning

$$\forall s, a \quad Q^*(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a')$$

↓

$$F(Q^*) = \frac{1}{|S| \cdot |A|} \sum_{s, a} \left(Q^*(s, a) - \left(r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^*(s', a') \right) \right)^2$$

$$s, a, r, s' \left(Q_{\text{new}}(s, a) = Q_{\text{old}}(s, a) - \lambda \left(Q_{\text{old}}(s, a) - \min_{a'} Q^*(s', a') \right) \right)$$

$$y(s, a) = r + \gamma \max_{a'} Q_{\text{old}}(s', a') - y(s, a) = (1 - \lambda) Q_{\text{old}}(s, a) + \lambda y(s, a)$$

Exploration - Exploitation Dilemma

$$a \sim \epsilon\text{-greedy}(Q(s, a)) : \left\{ \begin{array}{l} \text{random action with } \frac{\epsilon}{|A|} \\ \text{argmax}_a Q(s, a) \text{ with prob. } 1-\epsilon \end{array} \right.$$

$$a \sim \text{Softmax}\left(Q(s, a) / \tau\right) = \underbrace{\exp\left(\frac{Q(s, a)}{\tau}\right)}_{\sum_a \exp\left(\frac{Q(s, a)}{\tau}\right)}$$

$\tau \rightarrow +\infty$ Softmax \approx Uniform

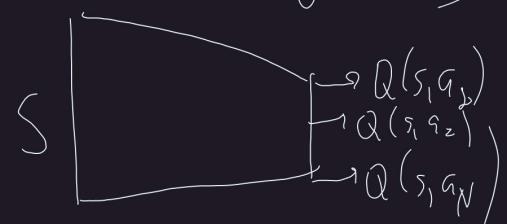
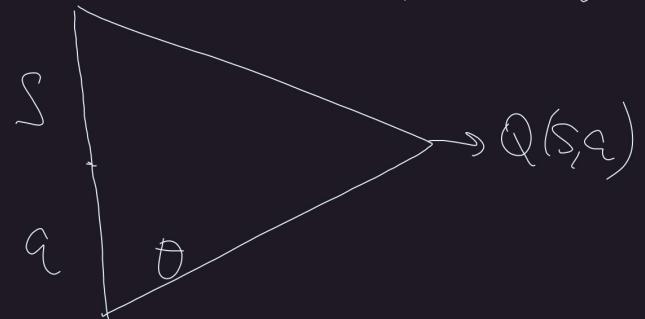
$\tau \rightarrow 0$ Softmax \approx argmax_a Q(s, a)

$$Q(s, a | \theta)$$

$$F(\theta) = \mathbb{E}_{s,a} \left[\frac{1}{2} \left(Q(s,a|\theta) - \left(r(s,a) + \gamma \mathbb{E}_{\substack{p(s'|s,a) \\ a'}} \max_{\theta'} Q(s',a'|\theta') \right) \right)^2 \right]$$

$$s, a, r, s'; \quad y(s,a) = r + \gamma \max_{a'} Q(s',a'|\theta) \quad \rightarrow \min_{\theta}$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \left(Q(s,a|\theta_{\text{old}}) - y(s,a) \right) \nabla_{\theta} Q(s,a|\theta)$$



Inference:

$$\pi(a|s) = \arg \max_a Q(s, a | \theta)$$

DQN (Deep Q Network)

Init. θ , buffer $M = \emptyset$ experience replay buffer

For episodes,

Init. s_0

For $t = 0, 1, 2, \dots, T$:

$$\bar{\theta} = \beta \bar{\theta} + (1 - \beta)\theta$$

$a_t \sim \epsilon\text{-greedy}(Q(s_t, a | \theta))$

get r_t, s_{t+1} ; $(s_t, a_t, r_t, s_{t+1}) \rightarrow M$

sample mini-batch $\{s_j, a_j, r_j, s'_j\}_{j=1}^B$ from M

for $j = 1, \dots, B$:

$$y_j = r_j + \gamma \max_{a'} Q(s'_j, a' | \bar{\theta})$$

$$\mathcal{L}_\theta = \mathcal{L}_\theta + \frac{1}{B} \sum_{j=1}^B (Q(s_j, a_j | \theta) - y_j)^2$$

$$\theta = \text{opt-step}(\nabla_{\theta} \mathcal{L})$$

$\bar{\theta}$ -target network

different from θ

Overestimation bias

