

# An Overview of Reinforcement Learning

Yue Zhao

November 27, 2019

## 1. Fundamental Knowledge

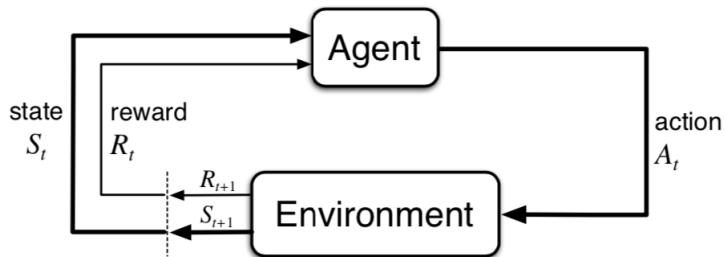
## 2. Value-based RL

- ▶ Monte Carlo
  - ▶ MCTS
  - ▶ \*AlphaGo Zero
- ▶ Time Difference
  - ▶ DQN

## 3. Policy-based RL

- ▶ Policy Gradient
- ▶ Trust Region Policy Optimization

# 1.1. Framework



## 1.2. Definitions

- ▶ Markov Decision Process:  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
- ▶ Policy  $\pi$ :
  - ▶  $\pi(a|s) = \mathbb{P}[A_t = a | S_t = s]$
- ▶ Return  $G_t$ :
  - ▶  $G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
- ▶ State-value function  $v_{\pi}(s)$ :
  - ▶  $v_{\pi}(s) = \mathbb{E}_{\pi}[G_t | S_t = s]$
- ▶ Action-value function  $q_{\pi}(s, a)$ :
  - ▶  $q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t | S_t = s, A_t = a]$

## 1.3. Properties

- ▶  $v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) q_{\pi}(s, a)$
- ▶  $q_{\pi}(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi}(s')$
- ▶  $v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$
- ▶  $q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$

## 1.4. Categorizing

- ▶ Value-based or Policy-based or AC
- ▶ On-policy or Off-policy
- ▶ Model-based or Model-free

## 2. Value-based

$$v_{\pi}(s) = E_{\pi} [G_t | S_t = s] = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

$$q_{\pi}(s, a) = E_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

## 2.1. Monte Carlo

$$\blacktriangleright v(s) = \frac{G_{11}(s) + G_{21}(s) + \dots}{N(s)}$$

$$\blacktriangleright v(s) = \frac{G_{11}(s) + G_{12}(s) + \dots + G_{21}(s) + \dots}{N(s)}$$



## 2.1. Monte Carlo

On-policy first-visit MC control (for  $\varepsilon$ -soft policies), estimates  $\pi \approx \pi_*$

Algorithm parameter: small  $\varepsilon > 0$

Initialize:

$\pi \leftarrow$  an arbitrary  $\varepsilon$ -soft policy

$Q(s, a) \in \mathbb{R}$  (arbitrarily), for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$  empty list, for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}(s)$

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots, 0$ :

$G \leftarrow G + R_{t+1}$

Unless the pair  $S_t, A_t$  appears in  $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$ :

Append  $G$  to  $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow \text{average}(Returns(S_t, A_t))$

$A^* \leftarrow \arg \max_a Q(S_t, a)$  (with ties broken arbitrarily)

For all  $a \in \mathcal{A}(S_t)$ :

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

evaluation

improvement

## 2.1.2. MCTS

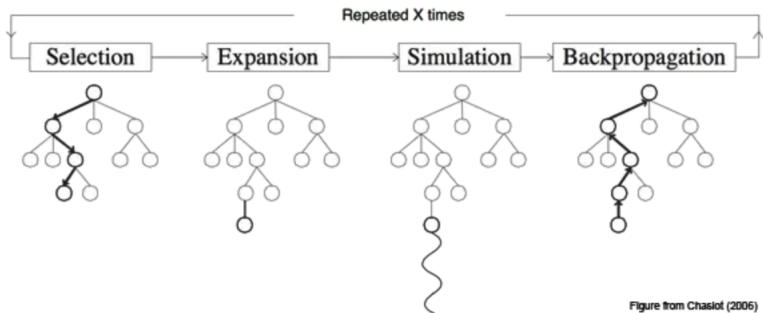


Figure from Chaslot (2006)

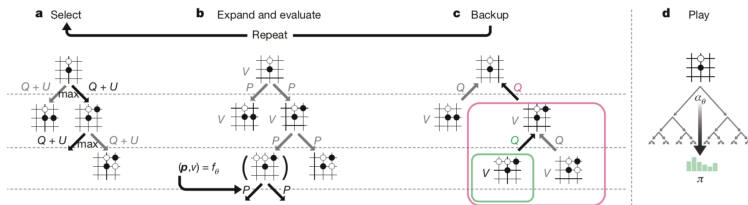
## 2.1.3. \*AlphaGo Zero

### ► NN:

$$(\mathbf{p}, v) = f_{\theta}(s)$$

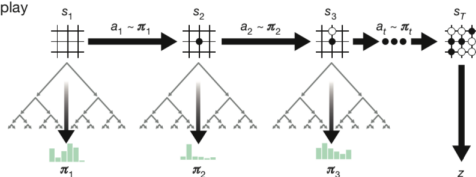
### ► Loss Function:

$$l = (z - v)^2 - \pi^T \log(\mathbf{p}) + c \|\theta\|^2$$

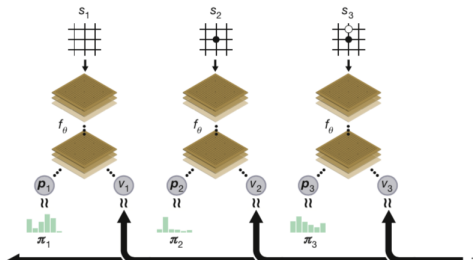


## 2.1.3. \*AlphaGo Zero

**a** Self-play



**b** Neural network training



## 2.2. Time Difference

- ▶ MC:  $V(S_t) \leftarrow V(S_t) + \alpha(G_t - V(S_t))$
- ▶ TD:  $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$

## 2.2. Q Learning

### Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

Algorithm parameters: step size  $\alpha \in (0, 1]$ , small  $\varepsilon > 0$

Initialize  $Q(s, a)$ , for all  $s \in \mathcal{S}^+$ ,  $a \in \mathcal{A}(s)$ , arbitrarily except that  $Q(\text{terminal}, \cdot) = 0$

Loop for each episode:

    Initialize  $S$

    Loop for each step of episode:

        Choose  $A$  from  $S$  using policy derived from  $Q$  (e.g.,  $\epsilon$ -greedy)

        Take action  $A$ , observe  $R, S'$

$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

$S \leftarrow S'$

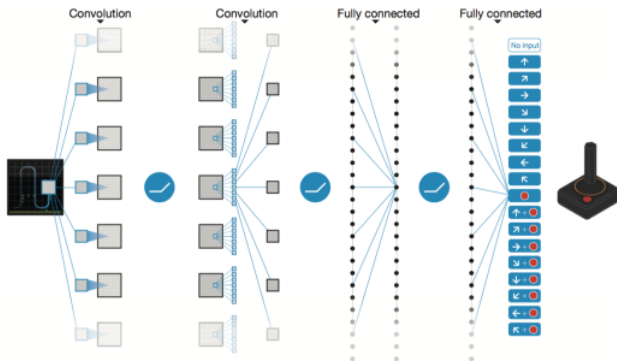
    until  $S$  is terminal

## 2.2. DQN

- ▶ Raw Pixel Input
- ▶ NN:  $Q(s, a) \approx f(s, a, w)$



## 2.2. DQN





## 2.2. DQN

- ▶ **Contributions:**

- ▶ Raw Pixel Input
- ▶ Experience Replay

## 2.2. DQN

$$\nabla_{\theta_i} L_i(\theta_i) = \mathbb{E}_{s, a \sim \rho(\cdot); s' \sim \mathcal{E}} [(r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i)]$$

---

**Algorithm 1** Deep Q-learning with Experience Replay

---

Initialize **replay memory**  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

**for** episode = 1,  $M$  **do**

    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$

**for**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$

        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$

        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3

**end for**

**end for**

---

### Improvements since Nature DQN

- ▶ **Double DQN:** Remove upward bias caused by  $\max_a Q(s, a, \mathbf{w})$ 
  - ▶ Current Q-network  $\mathbf{w}$  is used to **select** actions
  - ▶ Older Q-network  $\mathbf{w}^-$  is used to **evaluate** actions

$$l = \left( r + \gamma Q(s', \underset{a'}{\operatorname{argmax}} Q(s', a', \mathbf{w}), \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

- ▶ **Prioritised replay:** Weight experience according to surprise
  - ▶ Store experience in priority queue according to DQN error

$$\left| r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right|$$

- ▶ **Duelling network:** Split Q-network into two channels
  - ▶ Action-independent **value function**  $V(s, v)$
  - ▶ Action-dependent **advantage function**  $A(s, a, \mathbf{w})$

$$Q(s, a) = V(s, v) + A(s, a, \mathbf{w})$$

### 3. Policy-based RL

- **Core Concept:**  $\pi(a|s, \theta) = \Pr \{A_t = a | S_t = s, \theta_t = \theta\}$

## 3.1. Policy Gradient

- ▶ **Target:**  $\max J(\theta)$
- ▶ **Update:**  $\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$

## 3.1. Policy Gradient

### Policy Gradient Theorem

For any MDP, in either episodic cases or continuing cases,

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla_\theta \pi(a|s, \theta)$$

*Proof.*

## 3.1. Policy Gradient

```
function REINFORCE
  Initialise  $\theta$  arbitrarily
  for each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  do
    for  $t = 1$  to  $T - 1$  do
       $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$ 
    end for
  end for
  return  $\theta$ 
end function
```

## 3.1. Policy Gradient

### REINFORCE with Baseline

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (q_\pi(s, a) - b(s)) \nabla_\theta \pi(a|s, \theta)$$

- **Update:**  $\theta_{t+1} \doteq \theta_t + \alpha (G_t - b(S_t)) \frac{\nabla_\theta \pi(A_t|S_t, \theta_t)}{\pi(A_t|S_t, \theta_t)}$



## 3.2. Trust Region Policy Optimization

- ▶ **Key Point:** Monotonic Improvement Guarantee.

## 3.2. Trust Region Policy Optimization

**Lemma 1.** Given two policies  $\pi, \tilde{\pi}$ ,

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{\tau \sim \tilde{\pi}} \left[ \sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \quad (19)$$

This expectation is taken over trajectories  $\tau := (s_0, a_0, s_1, a_0, \dots)$ , and the notation  $\mathbb{E}_{\tau \sim \tilde{\pi}} [\dots]$  indicates that actions are sampled from  $\tilde{\pi}$  to generate  $\tau$ .



$$\eta(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A^{\pi}(s, a)$$



$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + E_{s \sim \rho_{old}, a \sim \pi_{old}} \left[ \frac{\tilde{\pi}_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)} A_{\theta_{old}}(s, a) \right]$$

## 3.2. Trust Region Policy Optimization

**Theorem 1.** *Let  $\alpha = D_{\text{TV}}^{\max}(\pi_{\text{old}}, \pi_{\text{new}})$ . Then the following bound holds:*

$$\eta(\pi_{\text{new}}) \geq L_{\pi_{\text{old}}}(\pi_{\text{new}}) - \frac{4\epsilon\gamma}{(1-\gamma)^2}\alpha^2$$

(8)

where  $\epsilon = \max_{s,a} |A_{\pi}(s, a)|$

## 3.2. Trust Region Policy Optimization

- ▶ maximize  $L_{\theta_{\text{old}}}(\theta)$   
subject to  $D_{\text{KL}}^{\text{max}}(\theta_{\text{old}}, \theta) \leq \delta$
- ▶ maximize  $L_{\theta_{\text{old}}}(\theta)$   
subject to  $\bar{D}_{\text{KL}}^{\rho_{\text{old}}}(\theta_{\text{old}}, \theta) \leq \delta$
- ▶ maximize  $\mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}, a \sim q} \left[ \frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{\text{old}}}(s, a) \right]$   
subject to  $\mathbb{E}_{s \sim \rho_{\theta_{\text{old}}}} [D_{\text{KL}}(\pi_{\theta_{\text{old}}}(\cdot|s) \parallel \pi_{\theta}(\cdot|s))] \leq \delta$

## 3.2. Trust Region Policy Optimization

$$\underset{\theta}{\text{maximize}} \left[ \nabla_{\theta} L_{\theta_{\text{old}}}(\theta) \Big|_{\theta=\theta_{\text{old}}} \cdot (\theta - \theta_{\text{old}}) \right] \quad (17)$$

$$\text{subject to } \frac{1}{2}(\theta_{\text{old}} - \theta)^T A(\theta_{\text{old}})(\theta_{\text{old}} - \theta) \leq \delta,$$

$$\text{where } A(\theta_{\text{old}})_{ij} =$$

$$\frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \mathbb{E}_{s \sim \rho_{\pi}} [D_{\text{KL}}(\pi(\cdot|s, \theta_{\text{old}}) \parallel \pi(\cdot|s, \theta))] \Big|_{\theta=\theta_{\text{old}}}.$$

$$\text{The update is } \theta_{\text{new}} = \theta_{\text{old}} + \frac{1}{\lambda} A(\theta_{\text{old}})^{-1} \nabla_{\theta} L(\theta) \Big|_{\theta=\theta_{\text{old}}},$$

## 3.2. Trust Region Policy Optimization

	<i>B. Rider</i>	<i>Breakout</i>	<i>Enduro</i>	<i>Pong</i>	<i>Q*bert</i>	<i>Seaquest</i>	<i>S. Invaders</i>
Random	354	1.2	0	-20.4	157	110	179
Human (Mnih et al., 2013)	7456	31.0	368	-3.0	18900	28010	3690
Deep Q Learning (Mnih et al., 2013)	4092	168.0	470	20.0	1952	1705	581
UCC-I (Guo et al., 2014)	5702	380	741	21	20025	2995	692
TRPO - single path	1425.2	10.8	534.6	20.9	1973.5	1908.6	568.4
TRPO - vine	859.5	34.2	430.8	20.9	7732.5	788.4	450.2