

Chapter 1. Why Reinforcement Learning

强化学习

Reinforcement Learning (RL) 根据背景做出最佳决策。它通过尝试自己学习。
它用于按次序决策，行动导致探索环境的情况，表现最佳。其优势在于
④ 它优化长期多步的奖励。它广泛用于优化模型，无需生成随机样本高
效拟合，因此十分迅速。

首个 RL 算法

Value estimate

$$V(s, w) = w_0 s_0 + w_1 s_1 + \dots + w_n s_n = w^T s, \quad \vec{s} \text{ 为输入向量，二进制}$$

Prediction error

$$\delta = \bar{E} - V(s, w)$$

\bar{E} 为实际结果

Weight update rule

$$w \leftarrow w - \alpha \delta x(s)$$

Chapter 2. Markov Decision Processes,

Dynamic Programming, and Monte Carlo Methods

Multi-Arm Bandit Testing

Policy Evaluation: 价值函数

行动 a 是集合 A 的元素，奖励 r 是 R 的元素，这些都是随机的

平均奖励计算: $r_{avg}(a) = \frac{1}{N(a)} \sum_{i=1}^{N(a)} r(a_i) = \frac{r_1 + \dots + r_{N(a)}}{N(a)}$

Inline value function:

$$r_{avg} \leftarrow \bar{r}_N \leftarrow \frac{1}{N} \left(r_{avg} \sum_{i=1}^{N-1} r_{i,0} \right) + r_{N-1}^{avg} + \frac{1}{N} (r_N - r_{N-1}^{avg})$$

更普遍来说，即 $r \leftarrow r + \alpha(r - r')$

Policy Improvement: 选最好的行动

ϵ -greedy bandit 算法

- 优势在
样本高
- 输入一个探索概率 ϵ ($\epsilon \in [0, 1]$)
 - 初始化 $r^{\text{avg}}(a) \leftarrow 0$, $N(a) \leftarrow 0$. 对每个 $a \in A$
 - 循环:
 - $x \leftarrow \begin{cases} \underset{a \in A(s)}{\operatorname{arg\max}} r^{\text{avg}}(as) & \text{概率 } 1 - \epsilon \\ \text{随机 } a & \text{概率 } \epsilon \end{cases}$
 - 执行 a , 获得奖励 r
 - $N(a) \leftarrow N(a) + 1$
 - $r^{\text{avg}}(a) \leftarrow r^{\text{avg}}(a) + \frac{1}{N(a)} (r - r^{\text{avg}}(a))$

贪婪模拟算法

function Environment($a, p(a)$)

输入 a , 概率 $p(a)$

输出 $r \sim 1 (p(a))$ 或 $r \sim 0 (1-p(a))$. 从概率

改进 ϵ -greedy 算法

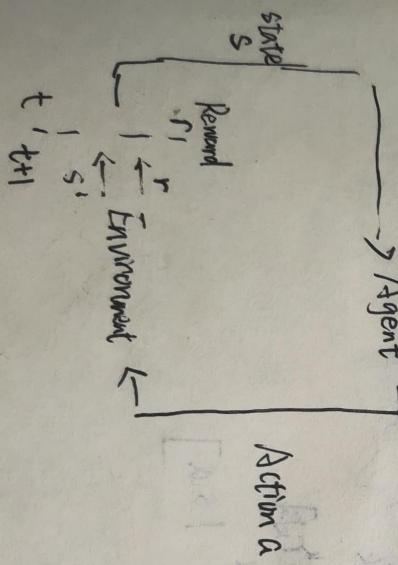
基于当前估计下的奖励分布选择行动更好, 因此以将概率分配行动代替直接选择行动. 以期望奖励作概率区分. 这里使用 softmax 函数

另一种是随时间减小 ϵ , 称为 annealing: 第三种探索未遇见过的情况 (VBS), 需要不需要使用 ϵ

马尔可夫决策过程 (MDP)

Agent 影响离散系统, 通过选择 action.
且行为选择一系列行动离散化系统

其中 state 和 reward 随机变量



MDP 转移模型: $p(s', r | s, a)$

实例：Inventory control

$$S = \{0, 1, 2\}, P_{restock} = \frac{1}{2}, A = \{\text{restock}, \text{none}\}$$

$$P(\text{sale}) \leftarrow 0.7$$

$$r = \begin{cases} 1 & S > 0 \text{ and a sale} \\ 0 & \text{otherwise} \end{cases}$$

MDP 矩阵表

$$\begin{array}{c|cc|cc|cc} & s' & p(s'|s,a) & r & & & \\ \hline s & a & & & & & \\ \hline 0 & none & 0 & 1-p & 0 & & \\ 0 & none & 0 & p & 0 & & \\ 0 & re & 1 & 1-p & 0 & & \\ 0 & re & 0 & p(sale) & 1 & & \\ \hline 1 & no & 1 & 1-p & 0 & & \\ 1 & no & 0 & p & 1 & & \\ \hline 1 & re & 2 & 1-p & 0 & & \\ 1 & re & 1 & p & 1 & & \\ \hline 2 & no & 2 & 1-p & 0 & & \\ 2 & no & 1 & p & 1 & & \\ \hline 2 & re & 2 & p & 1 & & \end{array}$$

转移矩阵

$$P(s', a=\text{none}) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & p & 1-p \\ 0 & 1-p & p \end{bmatrix}$$

$$P(s', a=\text{restock}) = \begin{bmatrix} p & 1-p & 0 \\ 0 & p & 1-p \\ 0 & 0 & 1 \end{bmatrix}$$

政策与价值函数

$$\text{回报 } G_t = r + r'_t + \dots + r'_T$$

$$\text{Discounted return } G_t = r + \gamma r'_t + \gamma^2 r''_t + \dots = \sum_{k=0}^{\infty} \gamma^k r_k$$

$$\text{State-Value 函数 } V_\pi(s) = \mathbb{E}_\pi[G_t | s] = \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_k \right]$$

$$\text{Action-Value 函数 } Q_\pi(s,a) = \mathbb{E}[G_t | s, a] = \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_k | s, a \right]$$

$$\text{Optimal value function } V^*(s) = \arg \max_a Q_\pi(s, a)$$

* 特卡罗政策生成

通过输入随机观察输出特征

On-policy MC 算法

1. input: a policy function, $\pi(a|s, \alpha(s, a))$
2. Initialize $\delta(s, a) \leftarrow 0$, for all $s \in \Phi, a \in A(s)$
3. loop:

(generate full episode trajectory following π)
 4. Initialize $G \leftarrow 0$
 5. loop for each step of episode, $t = T_1, T_2, \dots, t_0$
 $G \leftarrow V_{t+1}$
 if (s, a) not in $(s_0, a_0), \dots, (s_{t-1}, a_{t-1})$
 Append G to Returns (s, a)
 $\delta(s, a) \leftarrow \text{average}(\text{Returns}(s, a))$

Value Iteration

1. input, $P(s', r|s, a)$, stopping threshold $\theta > 0$
 2. initialize $V(s)$, for all $s \in \Phi$, $\nabla \leftarrow 0$
 3. do:
 4. loop for each $s \in \Phi$
 $v \leftarrow V(s)$
 $V(s) \leftarrow \max_a \sum_{s'} p(s', r|s, a)[r + \gamma V(s')]$
 $\nabla \leftarrow \max(\nabla, |V - V(s)|)$
- while $\nabla > \theta$
 output an optimal policy, $\pi(s) = \arg\max_{a \in A(s)} \sum_{s'} p(s', r|s, a)[\gamma V(s')]$

2 showing of transitions obtained

1. first 200 reward, 2000 2nd reward
 2. 1 reward bin - a little bit
 $[0.0, 0.2) \rightarrow [0.0, 0.2) \rightarrow [0.0, 0.2)$
 $[0.2, 0.4) \rightarrow [0.2, 0.4) \rightarrow [0.2, 0.4)$

labeled for 2nd

Chapter 3. Temporal Difference Learning,

π -Learning and n-Step algorithms.

SARSA

Temporal-Difference (TD) Learning 简介

Online ~~状态-动作~~ state-value ~~估计数~~

$$V_{\pi}(s) = \bar{E}_{\pi}[G(s)] \leftarrow E_{\pi}[V_{\pi}(s) + \alpha (r - V_{\pi}(s)) | s]$$

Online ~~动作-回报~~: state-value ~~估计数~~

$$V_{\pi}(s) = \bar{E}_{\pi}[G(s)] \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

TD state-value ~~估计数~~

$$V_{\pi}(s) = \bar{E}_{\pi}[V_{\pi}(s) + \alpha (r + \gamma V_{\pi}(s') - V_{\pi}(s)) | s]$$

Online TD state-value ~~估计数~~

$$V_{\pi}(s) \leftarrow V_{\pi}(s) + \alpha (r + \gamma \max_{a \in A(s)} V_{\pi}(s') - V_{\pi}(s))$$

(π -Learning)

它是 TD 的一个补充法

π -Learning online update rule to estimate the expected return

$$\delta(s, a) \leftarrow G(s, a) + \alpha (r + \gamma \arg\max_{a' \in A(s)} Q(s', a') - Q(s, a))$$

(π -Learning: off-policy TD)

1. Input: a policy that uses the action-value function, $\pi(a | s, Q(s, a))$

2. Initialize $Q(s, a) \leftarrow 0$, for all $s \in \mathcal{S}$, $a \in A(s)$

-loop: for each episode

Initialize environment to provide s

do:

Choose a from s using π , breaking ties randomly

Take action a , and observe r, s'

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \arg\max_{a' \in A(s')} Q(s', a') - Q(s, a)]$$

$$s \leftarrow s'$$

while s is not terminal

SARSA

提供更直接的 TD 估计

SARSA online estimate of the action-value function

$$Q(s, a) \leftarrow Q(s, a) + \alpha(r + \gamma Q(s', a') - Q(s, a))$$

SARSA (on-policy TD)

1. Input: a policy that uses the action-value function, $\pi(a|s, Q(s, a))$
2. Initialize $Q(s, a) \leftarrow 0$, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$
3. Loop:
 - a. Initialize environment to provide s
 - b. Choose a from s using π , breaking ties randomly
 - c. Take action a , and observe r, s'
 - d. Choose a' from s' using π , breaking ties randomly
 - e. $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma Q(s', a') - Q(s, a)]$
 - f. $s \leftarrow s', a \leftarrow a'$
- while s is not terminal

~~Following SARSA~~

n-step TD

2-step expected return

$$\tilde{G}_{t:t+2} = r + \gamma r' + \gamma^2 Q(s'', a'')$$

n-step expected return

$$G_t: t+n = r + \gamma r' + \gamma^2 r'' + \cdots + \gamma^{n-1} r_{n-1} + \gamma^n Q(s_n, a_n)$$

n-step TD update rule to estimate the expected return using the action-value function

$$Q(s, a) \leftarrow Q(s, a) + \alpha(G_{t:t+n} - Q(s, a))$$

NO

m.cn



考研数学
交流群

n -step SARSA

1. Input: a policy that uses the action-value function $\pi(a|s, Q(s,a))$
step size $0 < \alpha < 1$, a positive integer n
2. Initialize $Q(s,a) \leftarrow 0$, for all $s \in \mathcal{S}, a \in A(s)$
3. Loop for each episode:
Initialize $T \leftarrow \infty$, $t \leftarrow 0$, replay buffers for S and A , and initialize state s
Choose a from s using π , breaking ties randomly
do:
 if $t < T$:
 Take action a , and observe s' , r
 if s' is terminal
 $T \leftarrow t+1$
 else:
 Choose a' from s' using π , breaking ties randomly
 $I \leftarrow t-n+1$
 if $I+n < T$
 $G \leftarrow G + \gamma^n [G(s_{I:n}, a_{I:n})]$
$$Q(s_I, a_I) \leftarrow Q(s_I, a_I) + \alpha [G - Q(s_I, a_I)]$$

 $t \leftarrow t+1$, $s \leftarrow s'$, $a \leftarrow a'$
 while $T \neq T+1$

Eligibility Traces

SARSA(λ)

1. Input: a policy that uses the action-value function, $\pi(a|s, Q(s,a))$,
step size $0 < \alpha < 1$, a tracer decay rate $0 \leq \lambda \leq 1$
2. Initialize $Q(s,a) \leftarrow 0$, for all $s \in \mathcal{S}, a \in A(s)$
3. Loop for each episode
Initialize $Z(s,a) \leftarrow 0$, for all $s \in \mathcal{S}, a \in A(s)$, and default states
Choose a from s using π , breaking ties randomly

do:
 Take action a , and observe r, s'
 Choose a' from s' using π , breaking ties randomly
 $\delta \leftarrow r + \gamma Q(s', a') - Q(s, a)$
 $\hat{z}(s, a) \leftarrow z(s, a) + 1$
 for each $s \in \mathcal{S}$
 for each $a \in A(s)$
 $Q(s, a) \leftarrow Q(s, a) + \delta / \hat{z}(s, a)$
 $z(s, a) \leftarrow \gamma \hat{z}(s, a)$
 $s \leftarrow s'$, $a \leftarrow a'$
 while s is not terminal

Chapter 4 Deep π -Networks

Learning 或 SARSA 适用于小型、离散的状态，处理更大的最常用深度人工
 神经网络

Deep π -learning
 对 DQN 的改进
 Choosing an action with DQN

$$a \leftarrow \operatorname{argmax}_{a \in A(s)} Q(s, a; \theta)$$

DQN loss function

$$L(\theta) = \mathbb{E}[(y - Q(s, a; \theta))^2]$$

Random DQN

对 DQN 进行 6 项补充：其二为 double π -learning 和 nstep returns

最重要的是对 agents 重新引入概念

Distributional Bellman equation

$$Z(s, a) \leftarrow R(s, a) + \gamma Z(s', a')$$

Chapter 5. Policy Gradient Methods

直指学习 policy 的益处

"Optimal policies" 通过 policy 来指定每一个 action 的概率。Q-learning 通过张量
产生 optimal policy 来输出 policy。这一方法适用于 deterministic problems
但采用随机 policy 可通过随机观测进行随机选择。

另外，使用随机 policy 能干净地 encapsulate the exploration dilemma
更长远的好处是它的 action preferences 随时间平滑变化。

计算 policy 的梯度

在优化问题中，performance 用 $J(\theta)$ 表示， θ 为模型参数
policy objective function

$$J(\theta) = \mathbb{E}_{\pi} [G(s, a)]$$

Vanilla gradient ascent

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

Policy gradient 介绍

Gradient of a policy

$$\begin{aligned} \nabla J(\theta) &= \nabla \mathbb{E}_{\pi} [G(s, a)] \\ &= \nabla \sum_{a \in A} I(s, a) \nabla \pi(a|s) \\ &\propto \sum_{a \in A} [I(s, a) \nabla \pi(a|s)] \\ &\propto \sum_{a \in A} [a(s, a) \bar{\pi}(a|s) - \frac{\nabla \bar{\pi}(a|s)}{\bar{\pi}(a|s)}] \\ &\propto \sum_{a \in A} [\alpha(s, a) \nabla \pi(a|s) \nabla \ln \bar{\pi}(a|s)] \\ &\propto \mathbb{E}_{\pi} [\nabla G \nabla \ln \pi(a|s)] \end{aligned}$$

Policy gradient update rule

$$\theta \leftarrow \theta + \alpha \nabla \ln \pi(a|s, \theta)$$

policy functions

* 逻辑性 policy

$$\text{Logistic function} : f(x) = \frac{1}{1+e^{-x}}$$

或

$$\text{Logistic policy} \\ \pi_{\theta}(a|s) = \begin{cases} \pi_{\theta}(a=0|s) \\ \pi_{\theta}(a=1|s) \end{cases} = \begin{cases} \frac{1}{1+e^{\theta^T s}} \\ 1 - \frac{1}{1+e^{\theta^T s}} \end{cases}$$

Logistic policy gradient

$$\nabla \ln(\pi(a|s|\theta)) = \begin{bmatrix} s - s\bar{\pi}(\theta^T s) \\ -s\pi(\theta^T s) \end{bmatrix}$$

softmax policy

$$\pi_{\theta}(a|s) = \frac{e^{\theta^T s}}{\sum_a e^{\theta^T s}}$$

softmax policy gradient

$$\nabla \ln(\pi(a|s,\theta)) = s - \sum_a s\pi(\theta^T s)$$

Basic Implementations

(蒙特卡洛 Reinforce)

REINFORCE 算法

1. input: a differentiable parameterized policy $\pi(a|s,\theta)$, step size $\alpha > 0$

2. initialize θ arbitrarily
3. loop for each episode

generate trajectory following $\pi(\cdot|\cdot,\theta)$

for each step in episode ($t=0, 1, \dots, T$)

$$\text{choose } a_t \leftarrow \sum_{k=t+1}^{T-1} r_k$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \pi(a_t|s_t, \theta)$$

end for

Reinforce with Baseline

REINFORCE with baseline algorithm

1. input: a differentiable parameterized policy $\pi(a|s, \theta)$
a differentiable parameterized state-value function $V(s, w)$
step size $\alpha_w < \alpha_w < 0$ and $0 < \alpha_\theta < 1$

Eligibility

2. Initialize θ and w arbitrarily

3. loop for each episode:

Generate trajectory following $\pi(\cdot | \cdot, \theta)$:
for each step in episode $t \leftarrow 0, 1, \dots, T$:

$$G_t \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} r_k$$

$$\delta \leftarrow G_t - V(s, w)$$

$$w \leftarrow w + \alpha_w \delta \nabla V(s, w)$$

$$\theta \leftarrow \theta + \alpha_\theta \gamma^t \delta \nabla \pi(a_t | s_t, \theta)$$

end for

n -Step Actor-Critic and Advantage Actor-Critic (A2C)

n -step actor-critic

1. input: a differentiable parameterized policy $\pi(a|s, \theta)$
a differentiable parameterized state-value function $V(s, w)$

Step size α_w and $\alpha_\theta > 0$

2. Initialize θ and w

3. loop for each episode

Initialize $I \leftarrow 1$, $t \leftarrow 0$, and s

while s is not terminal

Choose a from s using π . breaking ties randomly

Take action a , observe s' , r

$$\delta \leftarrow (G_t : t \leftarrow n - V(s, \theta))$$

$$w \leftarrow w + \alpha_w \delta \nabla V(s, w)$$

$$\theta \leftarrow \theta + \alpha_\theta I \delta \nabla \ln \pi(a|s, \theta)$$

$$I \leftarrow \gamma I$$

$$s \leftarrow s'$$

Advantage function approximation for policy gradient algorithms

$$\begin{aligned} A^{\pi}(s, a) &= Q(s, a) - V(s) \\ &= r + Q(s', a') - V(s) \\ &= r + V_{\pi}(s') - V(s) \end{aligned}$$

Eligibility Traces Actor-Critic

Eligibility Actor-Critic with eligibility traces

- Input: a differentiable parameterized policy $\pi(a|s, \theta)$,
a differentiable parameterized state-value function $V(s, w)$
step size $\gamma < 1$, $\alpha_\theta < 1$, and trace-decay rates $0 < \lambda, \omega \leq 1$

2. initialize θ and w arbitrarily

3. loop for each episode

 Initialize $t \leftarrow 0$, $T \leftarrow 1$ and s

$z_w \in \mathbb{C}^J$. (same dimensionality as w) and

$z_\theta \in \mathbb{C}^J$. (same dimensionality as θ)

 while s is not terminal:

 choose a from s using π , breaking ties randomly

 Take action a , observe s' , r

$$\delta \leftarrow r + \gamma V(s', w) - V(s, \theta)$$

$$z_w \leftarrow \gamma \lambda_w z_w + \nabla V(s, w)$$

$$z_\theta \leftarrow \gamma \lambda_\theta z_\theta + I \nabla \ln \pi(a|s, \theta)$$

$$w \leftarrow w + \alpha_w \delta z_w$$

$$\theta \leftarrow \theta + \alpha_\theta I^\top \delta z_\theta$$

$I \leftarrow \gamma I$

$s \leftarrow s'$

Equivalent implements of the policy gradient algorithm

$$\nabla_{\theta} J(\theta) = E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) (\gamma - V_{\pi}(s))]$$

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) ((1 - \gamma)V_{\pi}(s) + \gamma)]$ REINFORCE

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) Q_{\pi}(s, a)]$ REINFORCE with Baseline

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) A_{\pi}(s, a)]$ Actor-Critic

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) \delta]$ Advantage Actor-Critic

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) \delta(t)]$ Temporal-difference Actor-Critic

= $E_{\pi_{\theta}} [\nabla_{\theta} \ln \pi_{\theta}(a|s) \delta_t]$ Eligibility Actor-Critic

Chapter 6. Beyond Policy Gradients

Off policy Algorithms

Importance sampling

$$E_f[x] = \sum_a x f(a)$$

$$E_g[x] = \sum_a x g(a)$$

$$= \sum_a \frac{g(a)}{f(a)} \cdot f(a)$$

$$= \bar{\theta} E_f \left[\frac{x g(a)}{f(a)} \right]$$

$$\approx \frac{1}{n} \sum_{i=1}^n \frac{g(x_i)}{f(x_i)}$$

Parameterized linear TD error

$$\delta = r + \gamma \theta^T \phi(s') - \theta^T \phi(s)$$

GTD(0) alg.

1. input: a policy that uses the parameterized action-value function.

$\pi(a|s, \phi; \theta)$. learning rate parameters. $0 < \alpha < 1$, $0 < \beta < 1$

2. initialize θ and u arbitrarily

3. loop: for each episode

 Initialize s

while s is not terminal:

choose a non s using π , breaking ties randomly

take action a , and observe r, s'

$$\delta \leftarrow r + \gamma \theta^\top \phi(s') - \theta^\top \phi(s)$$

$$q := w + \beta [\phi(s') - \theta^\top \phi(s)]$$

$$\theta \leftarrow \theta + \alpha (\phi(s) - q) \phi(s)^\top w$$

$$s \leftarrow s'$$

Greedy π^* algorithm

Greedy π^* algorithm

$$q := w + \beta [\phi(s') - \phi(s)]^\top w$$

$$q \leftarrow \theta + \alpha [\phi(s) - q] (w^\top \phi(s)) \phi(s)$$

$$s \leftarrow s'$$

Greedy Actor-Critics

$\hat{\pi}$ -policy gradient estimate

$$\nabla_{\theta} J(\pi_{\theta}) = E_{s \sim \pi_{\theta}, a \sim \pi} \left[\frac{\nabla_{\theta} \text{log} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \nabla_{\theta} \ln \bar{V}_{\theta}(a|s) \right]$$

Deterministic Policy Gradients

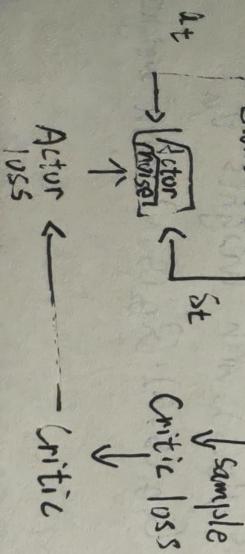
~~Stochastic~~ $\hat{\pi}$ -policy vanilla and deterministic actor-critic objective function gradients, for comparison

$$\nabla_{\theta} J(\pi_{\theta}) = E_{s \sim \pi_{\theta}, a \sim \pi} \left[\frac{\nabla_{\theta} \ln \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} V_{\theta}(\text{Intervals}) R_{\pi}(s, a) \right]$$

$$\nabla_{\theta} J(\pi_{\theta}) = E_{s \sim \pi_{\theta}} \left[\nabla_{\theta} \ln \pi_{\theta}(s, \mu_{\theta}(s)) \right]$$

$$= E_{s \sim \pi_{\theta}} [\nabla_{\theta} \mu_{\theta}(s) \nabla_{\theta} \ln \pi_{\theta}(s, a)]$$

DDPG π :



DDPG algorithm

1. input: an actor neural network, $\mu_\theta(s)$, with weights θ_μ ,
a critic neural network, $Q_\theta(s, a)$, with weights θ_α
a relay buffer B
the size of a training minibatch, N .
the network update rate, τ ,
the discount-rate parameter, γ ,
and a noise function for exploration, \mathcal{N} .
2. Initialize $\theta'_\mu, \theta'_\alpha, R$.
and the target weights, $\theta'_\mu \leftarrow \theta_\mu, \theta'_\alpha \leftarrow \theta_\alpha$
3. loop: for each episode

Initialize s

while s is not terminal:

$$a \leftarrow \mu_\theta(s) + \mathcal{N}$$

Take action, a , and observe r, s'

Store (s, a, r, s') in B

Sample a random minibatch of N transitions. (s_i, a_i, r_i, s'_i) from B

$$y_i \leftarrow r_i + \gamma Q_{\theta'_\alpha}(s'_i, \mu_{\theta'_\mu}(s'_i))$$

Update the critic by minimizing: $\frac{1}{N} \sum_i (y_i - Q_{\theta_\alpha}(s_i, a_i))^2$

Update the actor by minimizing: $\frac{1}{N} \sum_i (\nabla_{\theta_\mu} \mu_\theta(s_i, \mu_{\theta_\mu}(s_i)))^2$

$$\theta'_\alpha \leftarrow \tau \theta_\alpha + (1-\tau) \theta'_\alpha$$

$$\theta'_\mu \leftarrow \tau \theta_\mu + (1-\tau) \theta'_\mu$$

$$s \leftarrow s'$$

TD3 algorithm

1. input an actor neural network, $\mu_\theta(s)$, with weights θ_μ ,
two critic neural networks, $(Q_{\theta_1}(s, a), Q_{\theta_2}(s, a))$, with weights θ_1, θ_2

a relay buffer B ,

the size of a training minibatch N

the network update rate τ

the discount-rate parameter γ .

a noise clipping threshold c

and a noise injection or exploration N

2. Initialize $\theta_\mu, \theta_1, \theta_2$.

and the target weights $\theta'_\mu \leftarrow \theta_\mu, \theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2$

3. loop: for each episode

Initialize s

while s is not terminal

$$a \leftarrow \mu_\mu(s) + N$$

Take action a , and observe r, s'

store (s, a, r, s') in B

Sample a random minibatch of N transitions (s_i, a_i, r_i, s'_i) from B

$$\tilde{\pi}_i \leftarrow \mu_\mu'(s'_i) + \text{clip}(N, -c, c)$$

$$y_i \leftarrow r_i + \gamma^m \max_{j=1,2} Q_{\theta_j'}(s'_i, \tilde{\pi}_i)$$

Update both critics, $\theta_j \leftarrow 1, 2$, by minimizing: $\frac{1}{N} \sum_i (y_i - Q_{\theta_j}(s_i, a_i))^2$

If $+ \text{mod } 1$, then:

Update the actor by minimizing: $\frac{1}{N} \sum_i (\nabla_{\mu} \hat{\pi}_{\theta_\mu}(s_i, \mu_\mu(s_i)))^2$

Update both critics, $\theta_j \leftarrow 1, 2 : \theta_j' \leftarrow \tau \theta_j + (1-\tau) \theta_j'$

$$\theta'_\mu \leftarrow \tau \theta_\mu + (1-\tau) \theta'_\mu$$

end if

$$s \leftarrow s'$$

Trust Region Methods

Advantage actor-critic policy gradient

$$L_{\mu_\mu(j)} = E_{\pi_\theta} [\ln \pi_\theta(a|s) / \pi_\theta(s, a)]$$

Kullback-Leibler divergence

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \left(\frac{P(x)}{Q(x)} \right)$$

Reformulated advantage actor-critic loss function

$$\nabla J(\theta) = \bar{E}_{\pi_\theta} [\nabla h(\pi_\theta(s,a)) A(s,a)]$$

$$= \bar{E}_{\pi_\theta} \left[\frac{\nabla \pi_\theta(s,a)}{\pi_\theta(s,a)} A(s,a) \right]$$

Presenting A2C in terms of maximization of the objective

$$\underset{\theta}{\text{maximize}} \quad L^{\text{PPO}}(\theta) = \bar{E}_{\pi_\theta} \left[\frac{\pi_\theta(s,a)}{\pi_{\text{old}}(s,a)} A(s,a) \right]$$

$$\begin{aligned} & \text{subject to } \bar{E}_{\pi_\theta} [\text{KL}[\pi_{\text{old}}(\cdot|s) || \pi_\theta(\cdot|s)]] \leq \delta \\ & \text{Penalizing large changes in policy} \end{aligned}$$

$$\begin{aligned} & \text{Natural policy gradients optimization loss function} \\ & \underset{\theta}{\text{maximize}} \quad \bar{E}_{\pi_\theta} \left[\frac{\pi_\theta(s,a)}{\pi_{\text{old}}(s,a)} A(s,a) - \beta \text{KL}[\pi_{\text{old}}(\cdot|s) || \pi_\theta(\cdot|s)] \right] \\ & \text{Clipping large changes in policy according to the importance sampling ratio} \\ & L^{\text{CLIP}}(\theta) = E[\text{mcrc}(\theta) A(s,a), \text{clip}(\text{rc}(\theta), 1-\epsilon, 1+\epsilon) \hat{A}(s,a)] \end{aligned}$$

The final objective function for PPO

$$L^{\text{CLIP+VF+ts}}(\theta) = E [L^{\text{CLIP}}(\theta) + C_1 L^{\text{VF}}(\theta) + C_2 H[\pi_\theta](s)]$$

Asynchronous PPO algorithm

1. for iteration = 0, 1, ...
for agent = 0, 1, ..., N-1

Run policy π_θ in environment for T time steps

Compute advantage estimates A_0, \dots, A_{T-1}

Optimize L with respect to θ , for k epochs and minibatch size MN
 $\theta_{\text{old}} \leftarrow \theta$

-kita
Shan
P&E

Maxim
Me
Sort

Chapter 7 Learning All possible Policies with Entropy Methods

熵 Shannon entropy 表示随机信息量
随机变量的信息量

$$H(X) = -\sum_{x \in X} p(x) \log_b p(x)$$

Maximum Entropy Reinforcement Learning

Maximum entropy objective function

$$\pi^* = \arg \max_{\pi} \sum_{s,a} E[r + \gamma H(\pi(a|s))]$$

Soft Actor-Critic

Maximum entropy state-value function

$$V(s) = E[r + \gamma V(s') + \alpha H(\pi(a|s))]$$

$$= E[\alpha Q(s, a) + \alpha \log \pi(a|s)]$$

Maximum entropy action-value function

$$Q(s, a) = E[r + \gamma E[V(s')]]$$

Maximum entropy action-value objective function

$$J_\alpha(\theta) = E[\frac{1}{2} (\alpha \log(\pi_\theta(s, a)) - Q_\theta(s, a))^2]$$

$$= E[\frac{1}{2} (Q_\theta(s, a) - \alpha \log(r + \gamma E[V_\theta(s')]))^2]$$

ratio

SAC policy objective function

$$J_\pi(\phi) = E[\alpha \log(\pi_\phi(a|s)) - \beta H(s, a)]$$

SAC temperature parameter objective function

$$J_\alpha = E[-\alpha \log(\pi_\phi(s, a)) - \alpha H]$$

Extentions to Maximum Entropy Methods

Deep soft policy gradient

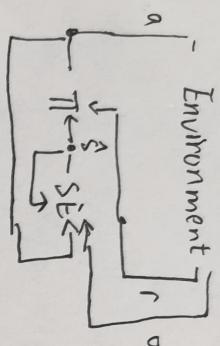
$$\nabla_\theta \phi = E[\nabla_\theta Q_\theta(s, a) - \log(\pi_\theta(s, a)) - 1] \nabla_\theta \log \pi_\theta(s, a)$$

NJ

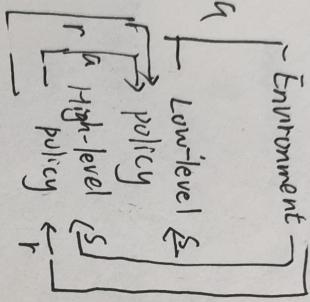
Chapter 8 Improving How An Agent Learns

Rethinking the MDP

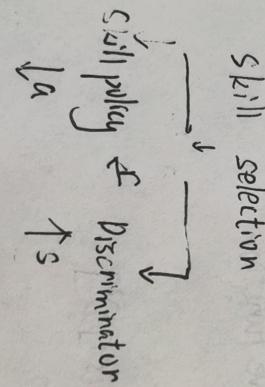
decompose a pMDP into an MDP and a state estimator



High-Low Hierarchies with Intrinsic Rewards (HLRU)



Schematic diagram of unsupervised skill learning framework (UPLAN)



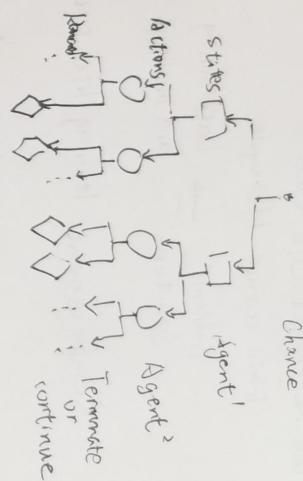
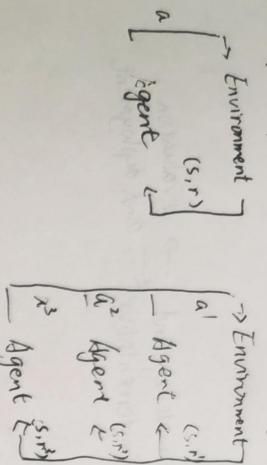
Environment

Schematic diagrams for different MARL frameworks.

(a) MDP

(b) Markov game

(c) Extensive-form game

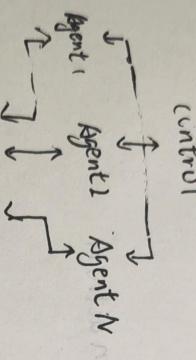


Centralized or Decentralized

A) Fully centralized

Learning

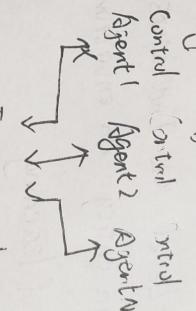
control



B) Centralized learning decentralized execution

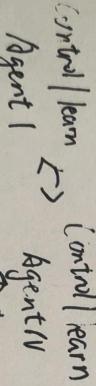
Learning

control

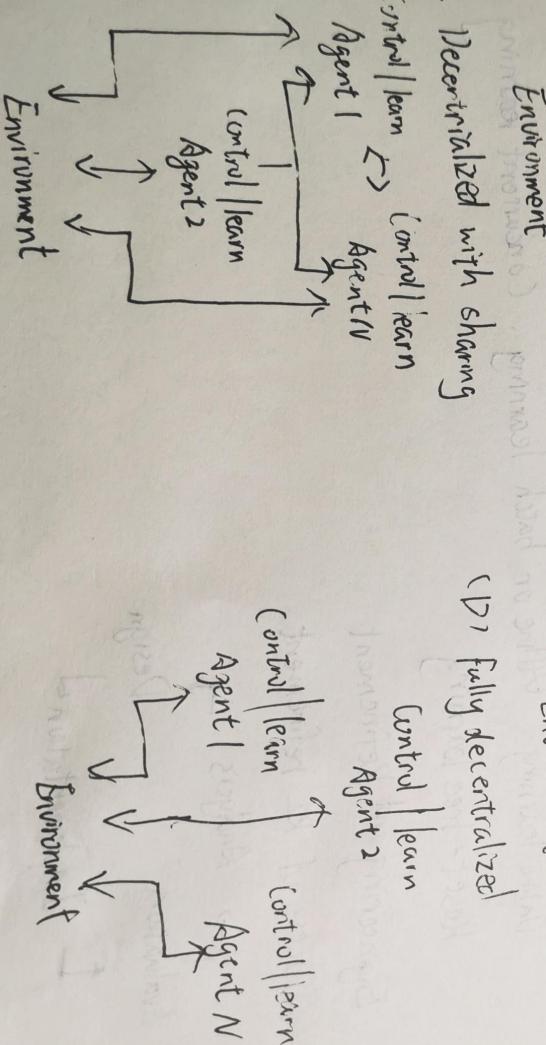
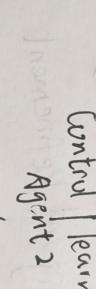


C) Decentralized with sharing

(control || learn)

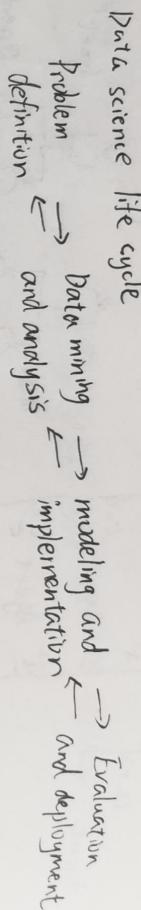


(control || learn)

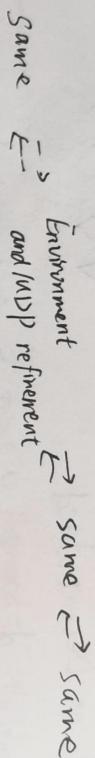


Chapter 9. Practical Reinforcement Learning

The RL Project Life Cycle



Reinforcement learning life cycle



RL Project : Problem definition

RL Problems ~~file~~ : Sequential , Strategic

Low-Level RL Indicators ,

An entity . An environment . A state . An action

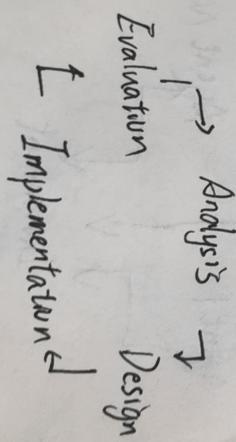
Types of learning :

Online learning . Offline or batch learning . Concurrent learning

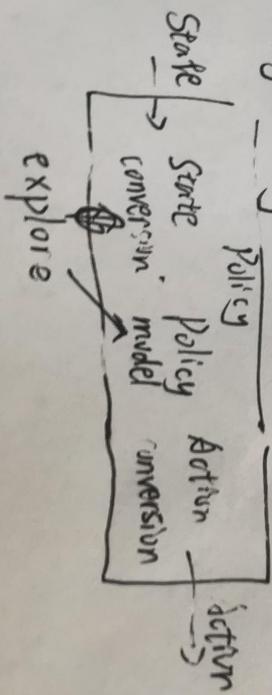
Reset-free learning

RL Engineering and Refinement

Process of RL refinement



Policy engineering



ning

nt

Implementation

Frameworks. Scenario . Evaluation

Chapter 10 Operational
Reinforcement Learning

Development

Tools , Architecture , Auxiliary Tooling Safety, Security and Ethics