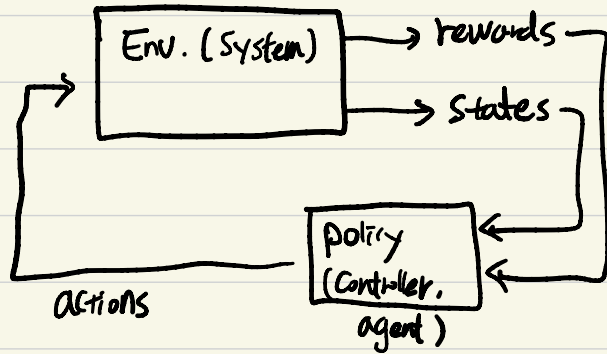


Reinforcement Learning

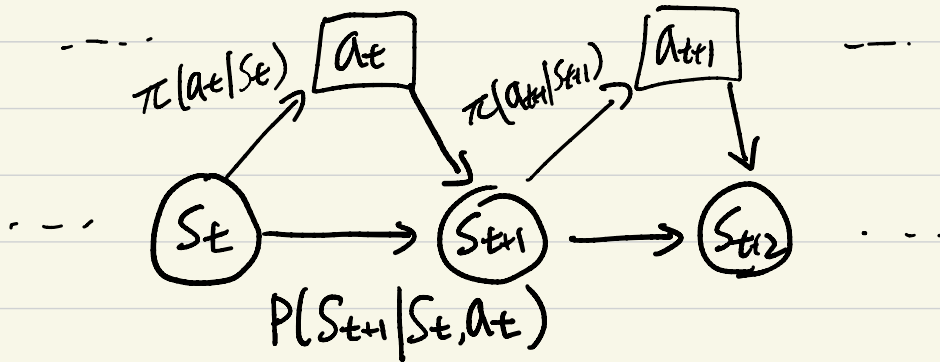
Today's

- MDP
- Q-Learning
- Deep Q-Learning



States, s_t
rewards, r_t
actions, a_t
transition prob. p
discount factor, γ

* Markov Decision Process (MDP)



Goal: train a policy to solve discrete MDP.

policy, π is a function that maps states to actions

optimal policy, π^*

① The return, R_t is the total discounted reward from time t .

$$R_t \doteq \sum_{k=0}^{T-1} \gamma^k r(S_{t+k}, A_{t+k})$$

② The state value function, $V^\pi(S_t)$ is the "expected" total discounted return from state S_t

$$\begin{aligned} V^\pi(S_t) &\doteq \mathbb{E}_\pi [R_t | S_t] \\ &= \mathbb{E}_\pi \left[\sum_{t'=t}^{T-1} r(S_{t'}, A_{t'}) | S_t \right] \\ &= \sum_{t'=t}^{T-1} \mathbb{E}_\pi [r(S_{t'}, A_{t'}) | S_t] \end{aligned}$$

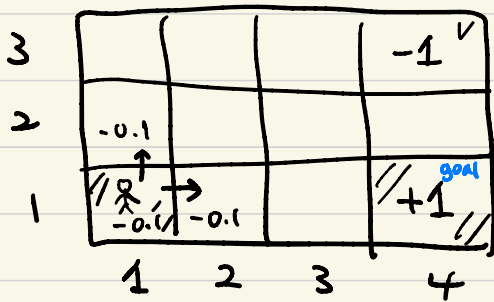
③ The state-action value function, $Q^\pi(S_t, A_t)$ is the expected total discounted return when action A_t is taken in state S_t .

$$\begin{aligned} Q^\pi(S_t, A_t) &\doteq \mathbb{E}_\pi [R_t | S_t, A_t] \\ &= \mathbb{E}_\pi \left[\sum_{t'=t}^{T-1} r(S_{t'}, A_{t'}) | S_t, A_t \right] \\ &= \sum_{t'=t}^{T-1} \mathbb{E}_\pi [r(S_{t'}, A_{t'}) | S_t, A_t] \end{aligned}$$

example,

$$\begin{aligned} Q(S_1, A_1) &= r(S_1, A_1) + \mathbb{E}_{S_2 \sim p(S_2 | S_1, A_1)} \left[\mathbb{E}_{A_2 \sim \pi(A_2 | S_2)} \left[r(S_2, A_2) | S_2 \right] \right] \\ &= r(S_1, A_1) + \mathbb{E}_{S_2 \sim p(S_2 | S_1, A_1)} V^\pi(S_2) \\ &\approx r(S_1, A_1) + V^\pi(S_2) \end{aligned}$$

Example: Grid World



- States: grid Coordinates $(1,1), (1,2) \dots (4,4)$
- Actions: $\{ \text{left, right, up, down} \}$
- Rewards: value in the box
- Trans. prob:

60%: you can go to next states what you intend.

20%: go right

20%: go left.

• objective: maximize R_t

• initial: $(1,1)$
States

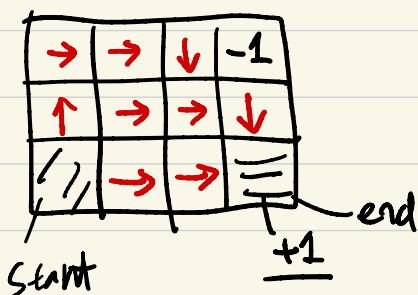
How do we find an optimal policy for Grid World.

Q-function tells you the value of starting in state S_t , applying action, A_t

$$Q^\pi(S_t, A_t) = \mathbb{E}_\pi \left[\sum_{t=0}^{T-1} r(S_t, A_t) \mid S_t, A_t \right]$$

$$= r(S_t, A_t) + \gamma \sum_{A_{t+1} \in \mathcal{A}} \pi(A_{t+1} | S_t) \sum_{S_{t+1} \in \mathcal{S}} P(S_{t+1} | S_t, A_t) V^\pi(S_{t+1})$$

The optimal Q-function for MDP, is deterministic



$$Q^*(S_t, A_t) = r(S_t, A_t) + \gamma \sum_{S_{t+1} \in \mathcal{S}} P(S_{t+1} | S_t, A_t) V^*(S_{t+1})$$

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

$$* Q^*(s,a) = r(s,a) + \gamma \sum_{s_{t+1} \in S} P(s_{t+1}|s_t, a_t) \max_{a_{t+1}} Q^*(s_{t+1}, a_{t+1})$$

What's the optimal action? a^* ?

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

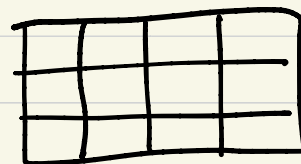
* Iterated Q-Learning

① Initialize V_0, Q_0

② For $i=1: \text{NUM_steps}$:

For $\forall s \in S$:

$$Q_i(s,a) \leftarrow r(s,a) + \gamma \sum_{s' \in S} P(s'|s,a) \max_{a'} Q_i(s',a')$$



Gridworld (2D)

once it's converged,

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

If we have a large number of states, then tabular representation is inefficient and intractable.

Quick Review

$$Q(s_t, a_t) = r(s_t, a_t) + \gamma \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1} | s_t, a_t) V(s_{t+1})$$

Standard Q-learning

$$Q(s, a) \leftarrow r(s, a) + \mathbb{E}_{s_{t+1} \sim P(s_{t+1} | s_t, a_t)} [V(s_{t+1})]$$



$$\leftarrow r(s, a) + \sum_{s_{t+1} \in \mathcal{S}} P(s_{t+1} | s_t, a_t) \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

How do we compute $\mathbb{E}_{s_{t+1} \sim P(s_{t+1} | s_t, a_t)} [V(s_{t+1})]$?

\Rightarrow The key idea in Q-learning is to use samples of the next state s_{t+1} in place of expectation.

$$\Leftrightarrow \mathbb{E}_{s_{t+1} \sim P(s_{t+1} | s_t, a_t)} [V(s_{t+1})] \approx \max_{a_{t+1}} Q(s_{t+1}, a_{t+1})$$

This will converge for the tabular case. Provable convergence of Q-learning requires visiting all state-action pairs infinitely many times.

* Q-learning w/ function Approximator.

• Q-learning

$$Q(S_t, a_t) \leftarrow r(S_t, a_t) + \max_{a_{t+1}} Q(S_{t+1}, a_{t+1})$$

- we parameterize Q function by func. approx, Q_θ .

- θ are the weights of the approximator,

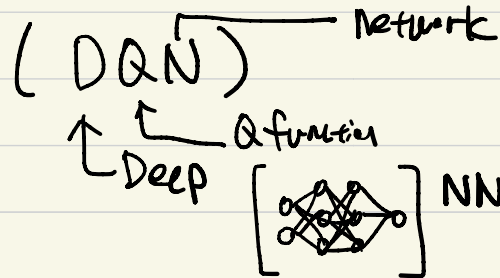
$$y_i \leftarrow r(S_{(i,t)}, a_{(i,t)}) + \max_{a_{(i,t+1)}} Q_\theta(S_{(i,t+1)}, a_{(i,t+1)})$$

$$\theta \leftarrow \underset{\theta}{\operatorname{argmin}} \sum_i \| y_i - Q_\theta(S_{(i,t)}, a_{(i,t)}) \|^2$$

- Use gradient descent to find θ .

• with this approach, the convergence is no longer guaranteed.

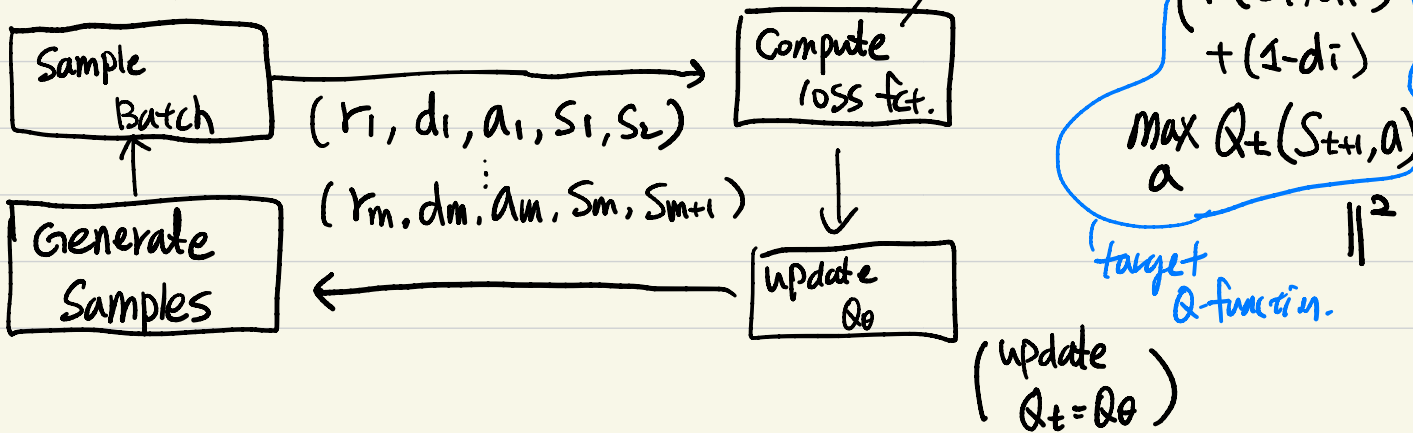
* Deep Q-learning



- ① Replay Buffers
- ② Target Q-function

$$(s, a, s', r) \quad \sum_i^m \| Q_\theta(s_i, a_i)$$

- Visual pseudo-code:



① Replay Buffer.

↳ save samples, $(S_i, A_i, S'_i, r_i, d_i)$

② Target Q-function

Idea: slow down the moving target value.

$$y_i \leftarrow r(S_i, a_i) + \max_{a'} Q_t(S'_i, a')$$

please check the class notes about DQN algorithm

Youtube link for DQN on Atari-games.

