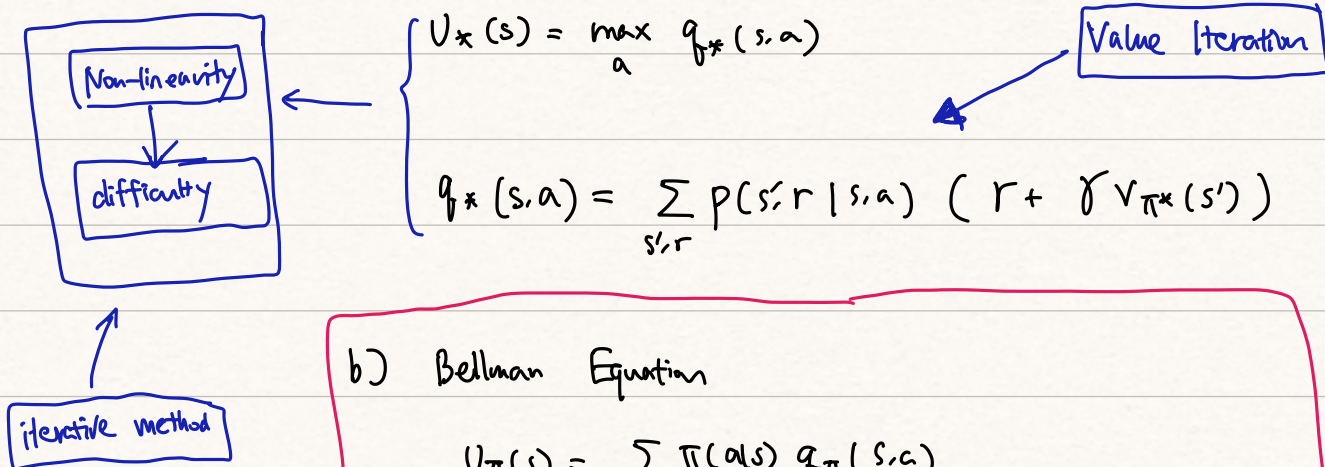


Model-based RL \Rightarrow know the knowledge of $p(s', r | s, a)$

Model-free RL $\Rightarrow p(s', r | s, a)$ is unknown

① Model-based RL \Rightarrow work with $\begin{cases} \text{Bellman Equa.} \\ \text{Bellman Optimality Equa.} \end{cases}$ exactly

Recap, a) Bellman Optimality Equation



b) Bellman Equation

$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$

$$q_{\pi}(s, a) = \sum_{s', r} p(s', r | s, a) (r + \gamma V_{\pi}(s'))$$

c) Optimal Policy

$$\pi^*(s) \in \operatorname{argmax}_{a \in A} q_*(s, a)$$

d) Policy Improvement

Policy Iteration

① Value Iteration

→ { Bellman Optimality Condition $V_*(s)$
Optimal Policy through $q_*(s, a)$

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



Matrix Form

$$\begin{cases} P(\pi)_{ss'} = p^\pi(s' | s) = p(s' | s, a) \\ b(\pi)_s = \mathbb{E}^\pi[r | s] = \sum_{s', r} p(s', r | s, a) \cdot r \end{cases}$$

$$a = \pi(s)$$

$$F(V) = \max_{\pi} [b(\pi) + \gamma P(\pi) \cdot V]$$

where π is deterministic policy

then V_* is optimal value function $\Leftrightarrow \underline{V_* = F(V_*)}$

$V_{k+1} = F(V_k)$ + Contraction Mapping Theorem



Convergence (Fixed Point)

② Policy Iteration

→ { Bellman Equation
Policy Improvement

Finite Algorithm! (for finite MDP)

$$\pi \xrightarrow{E} V_\pi \xrightarrow{I} \pi'$$

(E): $V_\pi(s) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) (r + \gamma V_\pi(s'))$



Matrix form

$$V_\pi = b(\pi) + \gamma P(\pi) V_\pi$$

linear Equation Solution



(I:)

$$\pi'(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a)$$

$$\pi' = \underset{\tilde{\pi} \in \text{deterministic}}{\operatorname{argmax}} [b(\tilde{\pi}) + \gamma P(\tilde{\pi}) \cdot V_{\pi}]$$

(policy improvement: $\pi' \geq \pi \Leftrightarrow V_{\pi'}(s) \geq V_{\pi}(s) \quad \forall s \in S$)

Model-free Algo

① Monte-Carlo Idea

→ estimate Expectation from Sample Mean

guarantee by $\begin{cases} \text{LLN} \\ \text{CLT} \end{cases}$

$$\underline{\mathbb{E}_{x \sim \mu} f(x) = \frac{1}{N} \sum_{i=1}^N f(X_i) \quad X_i \sim \mu \text{ (i.i.d.)}}$$

→ When $p(s', r | s, a)$ is our known knowledge,

then $V_{\pi} = b(\pi) + \gamma P(\pi) V_{\pi}$

→ When $p(s', r | s, a)$ is unknown,

then we go back to the defn of $V_{\pi}(s)$.

$$V_{\pi}(s) = \mathbb{E}_{\pi} [G_0 | S_0 = s]$$

⇒ Monte-Carlo Comes in!

↪ Sample N trajectories (episodes) through Policy π
 $\{S_t^{(n)}, R_t^{(n)}\} \quad n=1,2,\dots,N \quad t=1,2,\dots,T$

and $S_0^{(n)} \equiv s$

$$\Rightarrow \underline{V_{\pi}(s) \approx \frac{1}{N} \sum_{n=1}^N G_0^{(n)}} \rightarrow \text{Monte-Carlo Idea}$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{t=0}^{\infty} \gamma^t R_{t+1}^{(n)}$$

⇒ $q_{\pi}(s,a)$ can be estimated in the same fashion !!

Then, do Policy Improvement

$$\pi'(s) \in \arg\max_a q_{\pi}(s,a)$$

Remark: approximation of improved policy

GPI

since the q -function is approximated through Monte-Carl-!

Generalized Policy Iteration Framework

Not solve for V_{π} exactly

Short-term Monte Carlo Fashion

② Temporal-Difference Idea \leadsto not sample N complete trajectories
(Monte-Carlo)

iterate formula: $V(s) \leftarrow (1-\alpha)V(s) + \alpha(r + V(s'))$

where a sample from policy $\pi \rightarrow \pi(\cdot|s)$

and (s', r) is achieved by environment simulator

Observation: $V_\pi(s) = \mathbb{E}_\pi [G_0 | S_0 = s]$

$$= \mathbb{E}_\pi [r + \gamma G_1 | S_0 = s]$$

Interpretation:

Lemma:

$$a_{k+1} = (1-\alpha)a_k + \alpha b$$

$$0 \leq \alpha < 1$$

$$\Rightarrow \boxed{\lim_{k \rightarrow \infty} a_k = b} \quad (\text{converge exponentially})$$

idea: $V_{k+1} = (1-\alpha)V_k + \alpha \mathbb{E}_\pi [G_0 | S_0 = s]$

$$\Rightarrow \lim_{k \rightarrow \infty} V_k = \underline{\mathbb{E}_\pi [G_0 | S_0 = s]} \quad \text{exponentially}$$

General Idea

$\mathbb{E}_\pi [G_0 | S_0 = s] \approx G_0$ "one sample from environment simulator"

Illustration

$$= \mathbb{E}_\pi [r + \gamma G_1 | S_0 = s]$$

$$\approx r + \gamma \mathbb{E}_\pi [G_1 | S_0 = s] = r + \gamma \sum_{s'} P(\pi)_{ss'} \mathbb{E}_\pi [G_1 | S_1 = s']$$

$$\approx r + \gamma \mathbb{E}_\pi [G_1 | S_0 = s, S_1 = s']$$

$$= r + \gamma V_\pi(s')$$

one-term estimator

Q-learning Framework \leftarrow TD idea (off-policy)

$$q(s,a) \leftarrow (1-\alpha) q(s,a) + \alpha [r + \gamma \max_{a'} q(s',a')]$$

determine which (s,a) pair should be updated a'

$\left\{ \begin{array}{l} a \rightarrow \text{sample from } \pi(\cdot|s) \\ a' \rightarrow \text{sample from greedy policy} \end{array} \right.$

simulator (s',r)