

RL: Agent-oriented learning—learning by interacting with an environment to achieve a goal; Learning by trial and error, with only delayed evaluative feedback (reward).

- No supervisor, only a reward signal
- Feedback is delayed, not instantaneous
- Time matters, i.e., sequential decision making
- Agent's actions affect the subsequent data/feedback from the environment

RL Problem is a considerable abstraction of the problem of goal-directed learning from interaction with the environment. **RL methods/solutions** specify how the agent changes its policy as a result of its experience. The agent's **goal** is to maximize the total amount of reward it receives over the long run.

Bandit Problem (repeatedly take a choice, reward from stationary prob dist, maximize total reward)

Action-Value Methods

$$Q_t(a) = \frac{R_1 + R_2 + \dots + R_{K_a}}{K_a}$$

Policies: Greedy policy, ϵ -greedy policy. The greedy method performs significantly worse in the long run because it often gets stuck performing suboptimal actions. (Exploitation vs. exploration dilemma)

Softmax (choose action a at t -th play with prob:)

$$\frac{e^{Q_t(a)/\tau}}{\sum_{i=1}^n e^{Q_t(i)/\tau}}$$

τ is a positive parameter called the temperature

- Large temp: all (nearly) equiprobable selection
- Small temp: greedy action selection

Algorithm Implementation

$$\begin{aligned} Q_{k+1} &= \frac{1}{k} (R_k + kQ_k - Q_k) \\ &= Q_k + \frac{1}{k} [R_k - Q_k] \end{aligned}$$

Choose constant step size for non-stationary problem. Weight recent reward more heavily.

Markov Property: the future is independent of the past, given the present. A Markov Process is a seq of random states with Markov prop, a tuple $\langle S, P \rangle$.

State Transition Probability Matrix

$$\mathcal{P} = \begin{matrix} & \text{to} \\ \text{from} & \begin{bmatrix} \mathcal{P}_{11} & \dots & \mathcal{P}_{1n} \\ \vdots & & \\ \mathcal{P}_{n1} & \dots & \mathcal{P}_{nn} \end{bmatrix} \end{matrix}$$

Markov Reward Process (MRP): a Markov process with reward values, a tuple $\langle S, P, R, \gamma \rangle$

$$\mathcal{R}_s = \mathbb{E} [R_{t+1} | S_t = s]$$

Return: the total discounted reward from time t

$$G_t = R_{t+1} + \gamma R_{t+2} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- γ close to 0 leads to short-sighted evaluation
- γ close to 1 leads to far-sighted evaluation

State-Value Function: expected return starting from state s

$$v(s) = \mathbb{E} [G_t | S_t = s]$$

Bellman Equation for MRP

$$v(s) = \mathbb{E} [R_{t+1} + \gamma v(S_{t+1}) | S_t = s]$$

$$v(s) = \mathcal{R}_s + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'} v(s')$$

Bellman Equation in Matrix Form

$$v = \mathcal{R} + \gamma \mathcal{P} v$$

$$v = (I - \gamma \mathcal{P})^{-1} \mathcal{R}$$

Markov Decision Process (MDP): MRP with decisions, a tuple $\langle S, A, P, R, \gamma \rangle$

Policies: a distribution over actions given states

$$\pi(a | s) = \mathbb{P} [A_t = a | S_t = s]$$

$$\mathcal{P}_{s,s'}^{\pi} = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{P}_{ss'}^a$$

$$\mathcal{R}_s^{\pi} = \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{R}_s^a$$

Value Functions for MDP

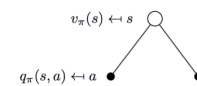
$$v_{\pi}(s) = \mathbb{E}_{\pi} [G_t | S_t = s]$$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a]$$

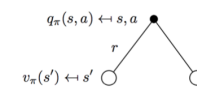
Bellman Expectation Equation for MDP

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s]$$

$$q_{\pi}(s, a) = \mathbb{E}_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) | S_t = s, A_t = a]$$



$$v_{\pi}(s) = \sum_{a \in \mathcal{A}} \pi(a | s) q_{\pi}(s, a)$$



$$q_{\pi}(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_{\pi}(s')$$

Bellman Expectation Equation in Matrix Form

$$v_{\pi} = \mathcal{R}^{\pi} + \gamma \mathcal{P}^{\pi} v_{\pi}$$

$$v_{\pi} = (I - \gamma \mathcal{P}^{\pi})^{-1} \mathcal{R}^{\pi}$$

Optimal Value Function

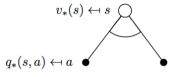
$$V_{*}(s) = \max_{\pi} V_{\pi}(s) \quad q_{*}(s, a) = \max_{\pi} q_{\pi}(s, a)$$

Optimal value specify the best performance in MDP

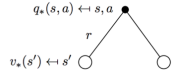
Optimal Policy

$$\pi_*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_{a \in \mathcal{A}} q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

Bellman Optimality Equation



$$V_*(s) = \max_a q_*(s, a)$$



$$q_*(s, a) = \mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_*(s')$$

Model-based RL

Dynamic Programming (DP) (Full Backup)

DP refers to a collection of algorithms that simplify a complicated problem by breaking it down into simpler sub-problems in a recursive manner.

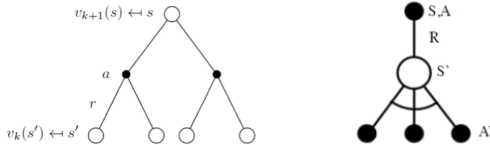
Policy Evaluation / Prediction

Problem: evaluate a given policy π

Solution: iterative application of Bellman

expectation backup $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_\pi$

Using synchronous backup: update $v_{k+1}(s)$ from $v_k(s')$



$$v_{k+1}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

Policy improvement

$\pi' = \text{greedy}(v_\pi)$

Policy Iteration / Control

Iterative policy evaluation + Greedy policy

improvement. Converge to v_*, π_* .

Generalized Policy Iteration (GPI)

Any policy evaluation + Any policy improvement.

Converge to v_*, π_* .

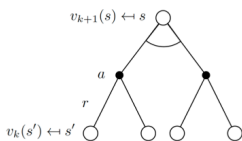
Value Iteration / Control

Problem: find optimal policy π

Solution: iterative application of Bellman

expectation backup $v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_*$

Using synchronous backup: update $v_{k+1}(s)$ from $v_k(s')$



$$v_{k+1}(s) = \max_{a \in \mathcal{A}} \left(\mathcal{R}_s^a + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}_{ss'}^a v_k(s') \right)$$

Asynchronous DP: (In-place DP; Prioritized sweeping; Real-time dynamic programming)

Model-free RL

Monte Carlo RL (MC)

MC Policy Evaluation / Prediction

Monte-Carlo policy evaluation uses empirical mean return instead of expected return.

First Visit MC Policy Evaluation, Every Visit MC Policy Evaluation. (in an episode)

$$N(s) \leftarrow N(s) + 1 \quad S(s) \leftarrow S(s) + G_t$$

$$V(s) = S(s)/N(s)$$

Incremental MC Updates: Update $V(s)$ incrementally after each episode. Useful in nonstationary problem

$$N(S_t) \leftarrow N(S_t) + 1$$

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} (G_t - V(S_t))$$

MC Policy Improvement / Control

Model Free GPI with MC Evaluation

Greedy policy improvement over $V(s)$ needs model, over $Q(s, a)$ is model free.

$$\pi'(s) = \operatorname{argmax}_{a \in \mathcal{A}} Q(s, a)$$

MC Policy Iteration

MC policy evaluation, $Q = q_\pi + \epsilon$ -greedy policy

improvement. Converge to q_*, π_*

Greedy in the Limit with Infinite Exploration (GLIE) Policy Iteration

All state-action pairs are explored **infinite** times.

Sample k-th episode: $\{S_1, A_1, R_2, \dots, S_T\} \sim \pi$

Evaluation:

$$N(S_t, A_t) \leftarrow N(S_t, A_t) + 1$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} (G_t - Q(S_t, A_t))$$

Improvement:

$$\epsilon \leftarrow 1/k$$

$$\pi \leftarrow \epsilon - \text{greedy}(Q)$$

Converge to $q_*(s, a)$

On policy: learn π from experience sampled from π

Off policy: learn π from experience sampled from μ

Temporal Difference (TD) (Sample Backup)

TD Prediction

$$V(S_t) \leftarrow V(S_t) + \alpha (R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$$

Learn online **every step**. Low variance, some bias.

$$\hat{\mathcal{R}}_{s,s'}^a = N(s, a, s')/N(s, a) \quad \hat{\mathcal{R}}_s^a = \sum r_s^a / N(s, a)$$

TD Control: Prediction + e-greedy

Sarsa Prediction

$$Q(S, A) \leftarrow Q(S, A) + \alpha (R + \gamma Q(S', A') - Q(S, A))$$

Q-Learning (Off-Policy) Prediction

$$Q(S, A) \leftarrow Q(S, A) + \alpha \left(R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$